

NEURAL DEHYDRATION: Effective Erasure of Black-box Watermarks from DNNs with Limited Data

Yifan Lu
Fudan University
Shanghai, China
luyf23@m.fudan.edu.cn

Wenxuan Li
Fudan University
Shanghai, China
liwx22@m.fudan.edu.cn

Mi Zhang*
Fudan University
Shanghai, China
mi_zhang@fudan.edu.cn

Xudong Pan
Fudan University
Shanghai, China
xdpan@fudan.edu.cn

Min Yang*
Fudan University
Shanghai, China
m_yang@fudan.edu.cn

ABSTRACT

To protect the intellectual property of well-trained deep neural networks (DNNs), black-box watermarks, which are embedded into the prediction behavior of DNN models on a set of specially-crafted samples and extracted from suspect models using only API access, have gained increasing popularity in both academy and industry. Watermark robustness is usually implemented against attackers who steal the protected model and obfuscate its parameters for watermark removal. However, current robustness evaluations are primarily performed under moderate attacks or unrealistic settings. Existing removal attacks could only crack a small subset of the mainstream black-box watermarks, and fall short in four key aspects: incomplete removal, reliance on prior knowledge of the watermark, performance degradation, and high dependency on data.

In this paper, we propose a watermark-agnostic removal attack called NEURAL DEHYDRATION (*abbrev.* DEHYDRA), which effectively erases all ten mainstream black-box watermarks from DNNs, with only limited or even no data dependence. In general, our attack pipeline exploits the internals of the protected model to recover and unlearn the watermark message. We further design target class detection and recovered sample splitting algorithms to reduce the utility loss and achieve data-free watermark removal on five of the watermarking schemes. We conduct comprehensive evaluation of DEHYDRA against ten mainstream black-box watermarks on three benchmark datasets and DNN architectures. Compared with existing removal attacks, DEHYDRA achieves strong removal effectiveness across all the covered watermarks, preserving at least 90% of the stolen model utility, under the data-limited settings, i.e., less than 2% of the training data or even data-free. Our work reveals the vulnerabilities of existing black-box DNN watermarks in realistic settings, highlighting the urgent need for more robust watermarking techniques. To facilitate

future studies, we open-source our code in the following repository: <https://github.com/LouisVann/Dehydra>.

CCS CONCEPTS

• Security and privacy → Security services; • Computing methodologies → Artificial intelligence.

KEYWORDS

Model Watermarking; Robustness; Removal Attack

ACM Reference Format:

Yifan Lu, Wenxuan Li, Mi Zhang, Xudong Pan, and Min Yang. 2024. NEURAL DEHYDRATION: Effective Erasure of Black-box Watermarks from DNNs with Limited Data. In *Proceedings of the 2024 ACM SIGSAC Conference on Computer and Communications Security (CCS '24)*, October 14–18, 2024, Salt Lake City, UT, USA. ACM, New York, NY, USA, 20 pages. <https://doi.org/10.1145/3658644.3690334>

1 INTRODUCTION

In recent years, deep neural networks (DNNs) are empowering real-world applications in computer vision [46, 64], natural language processing [16, 22], and autonomous driving [4, 5]. However, training a modern DNN from scratch requires time-consuming data collection and high-end computing resources. To protect the intellectual property of well-trained DNNs, *model watermarking* is an emerging tool for verifying the ownership of DNNs in case of model stealing or abusing.

Generally, in a DNN watermarking scheme, the secret watermark information is embedded into the *target model* at the watermark embedding stage during training. At the subsequent verification stage, the watermark can be extracted from the *suspect model* to determine its ownership. Depending on how the suspect model is accessed during verification, current watermark algorithms can be categorized into *white-box* and *black-box* watermarks. A white-box watermark is usually embedded into the suspect model's internal parameters [9, 54, 57] or neuron activation [15]. In comparison, a black-box watermark is embedded into a model's prediction behaviour on a set of specially-crafted samples (i.e., *watermark data*) by specifying their expected classification results (i.e., *target classes*) [2, 27, 34, 63]. Owing to the weaker access requirement for verification, black-box watermarks have gained increasing popularity in both academia and industry (e.g., IBM [19]), with various designs in watermark data and target class settings spurred.

*Corresponding authors: Mi Zhang and Min Yang.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).
CCS '24, October 14–18, 2024, Salt Lake City, UT, USA.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-0636-3/24/10
<https://doi.org/10.1145/3658644.3690334>

The robustness of black-box DNN watermarking schemes against *removal attacks*, which modifies the model parameters to cause verification failure, is crucial for their real-world trustworthiness. Unfortunately, current robustness evaluations are mostly performed under moderate attacks or unrealistic settings. Specifically, existing attacks are mainly categorized into three types: pruning-based [38, 54], finetuning-based [10, 47, 65] and unlearning-based [3]. However, these approaches could only crack a small subset of existing black-box watermarks (i.e., reducing the watermark verification success rate to almost random, with most of the utility preserved), which are unable to pose practical threats when the attacker is watermark-agnostic. To be more specific, these attacks may to varied degrees remove the watermark incompletely, require prior knowledge of the target watermark, hamper the model’s utility or have a strong dependence on the clean dataset. Recently, another type of model extraction-based attack has been demonstrated to effectively remove the watermarks when over 30% of the training data is available [42, 53]. However, as commercial DNNs are usually trained on private data, it is usually impractical to access such an amount of training data. A general comparison of their attack budgets is summarized in Table 1 and a more detailed analysis is presented in Section 2.2.

Our Work. In this paper, we propose a watermark-agnostic removal attack called NEURAL DEHYDRATION (*abbrev.* DEHYDRA), which effectively erases all ten mainstream black-box watermarks from DNNs, with only limited or even no data dependence. To achieve these objectives, we mainly address the following two challenges:

C1: How to design a watermark-agnostic effective removal attack?

Despite the complex designs of existing black-box watermarks in watermark data and target class settings, we identify a shared trait: they exploit the over-parameterization of DNNs to especially memorize the watermark data correlated to the target classes. Therefore, the model internals should be a sufficient source to reconstruct the watermark information. Inspired by this, we design a general-purpose attack framework to **recover and unlearn** the watermark data from the protected model, similar to the dehydration reaction in chemistry. Specifically, at the first stage, we use an aggressive model inversion technique to reverse-engineer class-wise samples close to the real watermark data from a target watermarked model (§4.1). At the second stage, we intentionally unlearn these samples during the finetuning process (§4.2). This basic DEHYDRA framework is effective against all the ten investigated black-box watermarks.

C2: How to preserve the model utility and reduce the data dependence during attack?

Despite effective, the basic attack might compromise the model’s utility, because the recovered samples could also contain information critical to the main task. To address this issue, we further enhance DEHYDRA in the following directions. ❶ We improve the recovering stage of DEHYDRA via incorporating **target class detection**. For watermarks with fixed target classes (such as [63] and [27], where all the watermark data are paired with the identical target class), we derive an observation called *watermark smoothness discrepancy* that, *the loss landscape of model is smoother on the target class than other classes*. We provide both theoretical and empirical analysis on the class-wise smoothness properties, and leverage the discrepancy to distinguish whether the target watermark has a fixed target class, and, if so, detect its target class. This allows us

Table 1: Comparison of attack budgets between our attack and existing watermark removal attacks, where ●/●/○ represent high, medium and low (or no) attack budget.

Attack Type	Removal Ineffectiveness	Watermark Knowledge	Utility Loss	Dataset Access
Prune	●	○	●	○
Finetune	●	○	●	●
Unlearn	○	●	○	●
DEHYDRA	○	○	○	○

to specially unlearn the recovered samples belonging to the target class when cracking a fixed-class watermark (§5.1). ❷ We improve the unlearning stage of DEHYDRA via **splitting the recovered samples**. Our design are based on the observed *normal data dominance* phenomenon that, *the class-wise recovered samples should be closer to normal samples compared with the watermark samples in that class*. Based on this, we develop a neuron-level criterion to carefully split the recovered samples into proxy watermark data and proxy normal data, with the proxy watermark data to ensure the removal effectiveness and the proxy normal data to preserve the original utility (§5.2). With the above improvements, DEHYDRA minimally impacts the model’s utility while ensuring its watermark-agnostic removal effectiveness. Moreover, the split proxy normal data can even allow for a data-free removal attack for watermarks identified with fixed target classes. Table 1 summarizes the advantages of our attack compared with existing removal approaches.

We conduct comprehensive evaluation of DEHYDRA against ten mainstream black-box watermarks on three benchmark datasets and DNN architectures, under three different data settings. Compared with six baseline removal attacks, our proposed DEHYDRA achieves strong removal effectiveness across all the covered watermarks, preserving at least 90% model utility, under the data-limited settings, i.e., less than 2% of the training data or even data-free.

To facilitate future studies, we open-source our code in <https://github.com/LouisVann/Dehydra>. Also, an extended version of this paper, with appendix included, is available in [40].

Our Contribution. We summarize the key contributions of this work as follows:

- We propose a watermark-agnostic removal framework, namely DEHYDRA, which exploits the model internals to recover and unlearn the underlying watermarks.
- Based on in-depth analysis, theoretical justifications and pilot studies, we further improve DEHYDRA with target class detection and recovered sample splitting algorithms, which help preserve the model utility and significantly reduce the dataset dependence.
- Extensive evaluations on various settings against ten mainstream black-box watermarks show that DEHYDRA is generally effective and has minimal utility impact, with much more relaxed or even no assumptions on the dataset availability.

2 BACKGROUND

2.1 DNN Watermarks

Traditional model protection techniques employ proactive strategies to prevent unauthorized access, such as encrypting model parameters or restricting access to APIs. However, these methods

have been proven by subsequent research to face leakage risks [50, 53]. To further claim and trace the ownership of DNN models, model watermarking technology has emerged.

A white-box watermarking scheme typically embeds a sequence of secret messages into the parameters or neural activations of the target model, and thus requires white-box access of the suspect model to extract the watermark [9, 15, 54, 57]. On the contrary, a black-box watermark (sometimes also called backdoor-based [39, 47], poisoning-based [21] or trigger set-based [34] watermarks) is embedded in the model's prediction behavior on a set of specially-crafted secret samples (i.e., *watermark data*) by specifying their expected prediction results (i.e., *target classes*) [2, 27, 30, 63]. Therefore, the verification process of black-box watermarking only requires the access to the prediction API.

To build a trustworthy ownership verification process, an ideal watermarking scheme should satisfy at least the following three key requirements [42, 59]. **① Fidelity:** The performance degradation on the target model should be as low as possible after the watermark is embedded. **② Integrity:** Models which are trained independently without access to the target model should not be verified as containing the watermark. **③ Robustness:** The embedded watermark in the target model should be resistant to potential removal attacks. In the worst cases when the watermark is removed, the utility of the model should decrease dramatically.

We mainly study on the black-box watermarks due to their more realistic settings. Existing black-box watermark schemes have diverse designs in the choices of the *watermark data* and the *target label*. In terms of watermark data, some work leverages specially-crafted patterns [27, 63] or random noise [63] as watermark patterns. Other black-box watermarks may use out-of-distribution samples as watermark data [2, 63]. Besides, in-distribution clean samples [45], adversarial samples near the decision boundary [30] or images with an imperceptible logo embedded by an encoder [37] can also be exploited as watermark data. In terms of the target class setting, existing black-box DNN watermarks can be categorized into fixed-class watermarks, where all the watermark data are paired with the identical target class (such as [27, 35, 63]), and non-fixed-class watermarks, where each watermark data is paired with its own target class (such as [2, 20, 30, 37, 45]). These complex designs pose additional challenges to a practical attacker, who usually has no knowledge of the underlying watermark schemes.

2.2 Watermark Removal Attacks

Black-box model watermarks are to some extent similar to backdoor attacks, as they both establish connections between some specified data and the target labels. Therefore, classical methods in backdoor defense, such as pruning, finetuning or trigger reverse-engineering, are usually considered potential threats for black-box watermarks [2, 27, 63]. Recently, some attacks especially targeting black-box watermarks are also proposed.

Existing watermark removal attacks can be mainly categorized into three types as follows, with the pros-and-cons of attacks from each category summarized in Table 1.

- **Pruning-based Attacks:** They prune the redundant weights [54] or neurons [38] in DNNs. However, to invalidate the underlying

watermarks, these methods need to prune a large proportion of weights or neurons, which causes an unacceptable utility loss.

- **Finetuning-based Attacks:** They continue to train the target model for a few epochs with some carefully designed finetuning techniques, such as learning rate schedule [10, 11], dataset augmentation [39], weight regularization [47], or continue learning with attention distraction [65]. Unfortunately, these attacks could only breach a subset of watermarks. Furthermore, due to stringent training restrictions, they require a large amount of clean data to maintain utility.
- **Unlearning-based Attacks:** They are mainly designed for black-box watermarks whose watermark data have the fixed pattern and the identical target class [3, 49]. For instance, the *Laundering* attack [3] uses the classical trigger reverse-engineering technique in backdoor defense literature *NeuralCleanse* [56], followed by an unlearning process. However, these attacks pose strong assumptions on the forms of underlying watermarks and are therefore not applicable to physical scenarios where attackers have no knowledge of the target watermarks. Also, current unlearning-based attacks all need access to source training data, which further limits their applicable scenarios.

Note that there also exists another type of attack aiming at training a new surrogate model using the knowledge transferred from the given watermarked model, i.e., **model extraction attacks** [27, 42, 60]. However, extraction attacks typically require a large query dataset and entail heavy computation costs [47]. Therefore, we do not include these attacks in the main experiment, but will make a comparison with them under different data settings in Section 7.

3 SECURITY SETTINGS

3.1 Mechanism of Black-box Watermarking

The true model owner is denoted as \mathcal{O} . During the training, \mathcal{O} embeds a watermark into the target model by training on both the clean dataset $X = \{(x_i, y_i)\}_{i=1}^{N_1}$ ($y_i \in \{0, 1, \dots, C-1\}$) from its main task, where C is the total number of the classification classes, and on a set of specially-crafted watermark data $X_w = \{(x_{w_i}, y_{w_i})\}_{i=1}^{N_2}$ ($y_{w_i} \in \{0, 1, \dots, C-1\}$), where target classes y_{w_i} can be either identical or sample-specific. We denote the watermarked model as f_w . In face of potential copyright infringement, \mathcal{O} hopes to verify the model ownership by comparing the suspect model's predictions on X_w with the pre-defined target labels, via the provided API (i.e., black-box access), as shown in Figure 1.

By inspecting the official watermark implementations, we observe the following technical designs which the owner \mathcal{O} will take to meet the three requirements in Section 2.1. Specifically, **①** to satisfy the **fidelity** purpose, \mathcal{O} usually choose an over-parameterized model with enough capacity, which is trained on both normal samples and watermark samples, with normal ones being dominant in quantity [2, 37], i.e., $N_2 \ll N_1$. **②** To ensure the watermark **integrity**, \mathcal{O} typically crafts watermark samples (x_{w_i}, y_{w_i}) that are significantly different from the normal samples in the target class y_{w_i} , whose size N_2 cannot be too small as well, to claim the ownership with high confidence [27, 35, 42]. **③** To enhance the **robustness** of the embedded watermark, \mathcal{O} tends to intentionally increase the involvement of the watermark data in the training

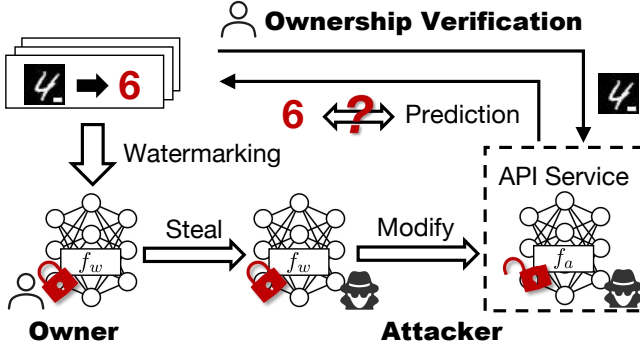


Figure 1: An illustration of our threat model. The model owner trains a watermarked model f_w with a secret set of watermark samples, and hopes to perform ownership verification via API queries for suspect models. On the other hand, the attacker has managed to acquire the white-box f_w , and aims to derive a surrogate model f_a with the watermark removed, evading subsequent ownership verification.

process. For instance, both Adi et al. [2] and Lederer et al. [33] concatenate the independently sampled normal data and watermark data into a data batch during each training iteration. Jia et al. [27] intentionally train the model on watermark samples and on normal samples alternatively. Lukas et al. [42] explicitly boost the ratio of watermark data when constructing the total training set [41].

3.2 Threat Model

Attack Scenario. In Figure 1, the attacker \mathcal{A} is unwilling to train a model of their own, due to the lack of training data, computation resource or expertise. Instead, \mathcal{A} attempts to infringe on the intellectual property of f_w trained by \mathcal{O} , and then host a similar service illegally or abuse it to make a profit, while evading the subsequent watermark verification process.

Similar to the settings in many watermark robustness studies [3, 10, 34, 42, 65], we assume the attacker has managed to acquire the white-box access to the target model f_w , i.e., they can observe the model parameters as well as modify them. This access could be achieved by, e.g., directly downloading the model after dishonestly agreeing with the official license (on open-source model hubs such as Hugging Face [17], DNNs might be released under licenses permitting research but prohibiting commercial use [14]), decrypting the possibly encrypted model from memory (DNNs on end devices) [50], or stealing functional replicas of DNNs via API access (MLaaS) [53]. Additionally, we assume that \mathcal{A} is aware of the existence of the watermark in f_w , and therefore their goal is to derive a surrogate model f_a from f_w with similar performance while invalidating the ownership verification.

Attack Budget. Specifically, as shown in Table 1, ❶ \mathcal{A} has no knowledge of the specific watermark algorithm or the watermark data X_w that \mathcal{O} used (i.e., watermark-agnostic). ❷ \mathcal{A} hopes to significantly lower f_a 's watermark verification accuracy on X_w , ❸ while preserving its clean accuracy on the main task. To make this attack practical, ❹ \mathcal{A} also wants to relax the dataset access requirement as much as possible.

Previous removal attacks may require a substantial subset of the clean training data X [2, 10, 63], or a considerable size of unlabeled data collected from open sources as proxy data [21, 47]. However, they usually ignore the difficulty of collecting high-quality data at such a scale. In this work, we mainly consider the attacker with limited data access, as in [3, 65]. The following three assumptions on dataset availability characterize attackers from being limited in knowledge to almost of zero knowledge:

- **In-distribution:** A limited number of correctly labeled samples (either manually by humans, or automatically by the well-training watermarked model) are available from the same distribution of the target model's main task. This setting is realistic since these in-distribution data are often hard to acquire.
- **Transfer:** A small amount of unlabeled data from a transfer distribution is available. Note that the distribution of these unlabeled data also matters, since finetuning the target model on samples from a completely different distribution might lead to catastrophic forgetting [28].
- **Data-free:** No additional data is available. This setting further relaxes the dataset access requirement for the attacker. To the best of our knowledge, our work is the first to formally consider data-free black-box watermark removal settings.

4 GENERAL FRAMEWORK OF DEHYDRA

In general, our DEHYDRA adopts the following two-stage attack pipeline: (1) **Watermark Recovering:** We first leverage model inversion technique to recover some samples close to real watermark data (§4.1). (2) **Watermark Unlearning:** We next deliberately unlearn these samples during the finetuning process, along with the auxiliary dataset \mathcal{A} possesses (§4.2). The procedures bear a resemblance to dehydration reactions in chemistry, where water molecules are expelled from compounds, inducing the formation of new substances. An overview of our general attack framework is shown in the upper region of Figure 2.

4.1 Watermark Recovering

Despite the complexity of various designs of watermarks, including different watermark patterns and strategies for setting the target classes, we summarize their commonality that, to implement robustness, existing watermarks exploit the over-parameterization property of DNNs to intentionally memorize the watermark data correlated to the pre-defined target labels. Note that training samples with higher repetition and memorization levels may face greater reconstruction risks [7]. Inspired by this, we believe the model internals should be a sufficient source to reconstruct the watermark information, and hence attempt to use model inversion techniques to reverse-engineer samples close to the real watermark data from a watermarked model.

Concretely, we use optimization-based model inversion to recover the samples hidden in the watermarked model f_w . Given a target class c and a batch of randomly initialized samples $B = \{\hat{x}_i\}_{i=0}^{M-1}$ (\hat{x}_i keeps the same shape as the images of f_w 's main task, while M is the number of samples to be optimized for the c -th class, we try to reverse-engineer samples which are similar to the watermark

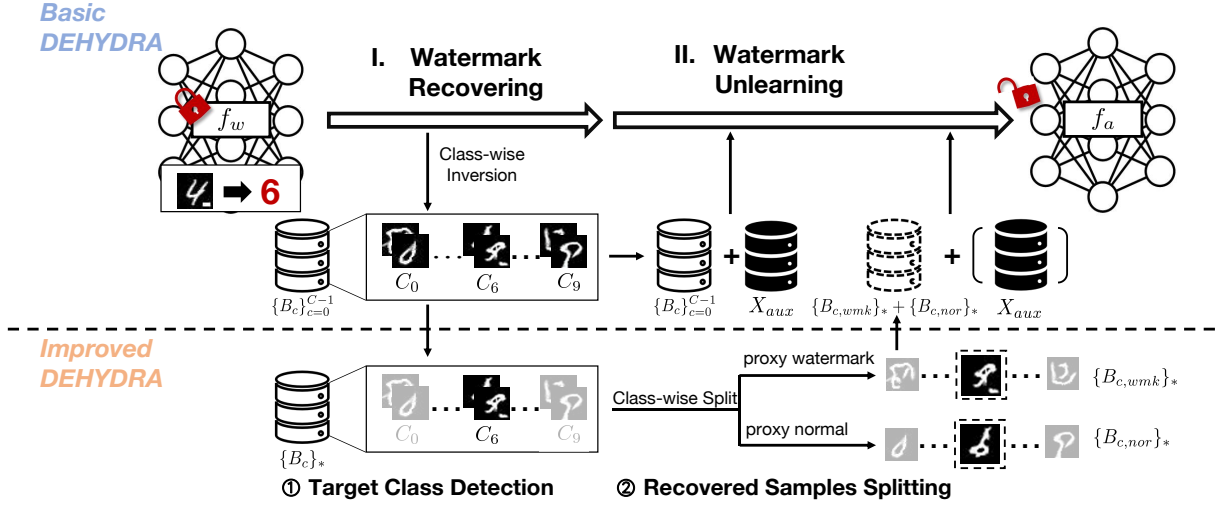


Figure 2: The overview of our DEHYDRA. The upper region shows our basic attack, comprising two stages, watermark recovering (§4.1), which reconstructs batch samples $\{B_c\}_{c=0}^{C-1}$ close to real watermark data from each class, and watermark unlearning (§4.2), which unlearns the recovered samples during finetuning, along with the auxiliary dataset X_{aux} . The lower region shows our improved designs, target class detection (§5.1) and recovered samples splitting (§5.2). The former detects the target classes after recovering (e.g., C_6 in this case), and only retains batches $\{B_c\}_*$ from those classes. The latter performs class-wise splitting on each B_c before unlearning, to identify samples closer to watermark or normal data, i.e., $\{B_{c,wmk}\}_*$ and $\{B_{c,nor}\}_*$, respectively.

data correlated to the target class c by optimizing

$$\min_B \left(\sum_{\hat{x}_i \in B} \mathcal{L}(f_w(\hat{x}_i), c) \right) + \mathcal{R}(B), \quad (1)$$

where $\mathcal{L}(\cdot)$ is the cross-entropy loss, and the regularizer $\mathcal{R}(\cdot)$ poses our prior knowledge on the watermark data.

In addition to classic priors on natural images (e.g., ℓ_2 regularization and total variation regularization [44]), our recovering algorithm mainly exploits the information stored in the batch normalization (BN) layers [26], as is done in [61, 62]. BN layers are widely adopted in mainstream DNNs [23, 24, 51], and keep useful statistics of the training data, capturing both low-level and high-level features in the neural network. Since the owner \mathcal{O} tends to intentionally increase the involvement of the watermark data during the training, the watermark data X_w are expected to play an important role in the BN statistics of the target model f_w .

To better utilize the watermark information hidden in the BN layers, we employ another feature regularization term $\mathcal{R}_{bn}(\cdot)$ to guide the recovered samples closer to the real watermark data:

$$\mathcal{R}_{bn}(B) = \sum_{l=1}^{L-2} \|\mu_l(B) - \mu_l\|_2^2 + \|\sigma_l^2(B) - \sigma_l^2\|_2^2, \quad (2)$$

where L is the total number of BN layers in f_w , μ_l and σ_l^2 are the running batch-wise mean and variance stored in the l -th BN layer, and $\mu_l(\cdot)$ and $\sigma_l^2(\cdot)$ are current statistics of feature maps before the l -th BN layer. According to [25, 48], the last several layers of a DNN capture high-level features, which are more correlated with the model's prediction behavior. To encourage the exploration in class c 's decision space and capture more watermark information, we intentionally dismiss the guidance of the last two BN layers.

In summary, the regularization term in Equation 1 is finally:

$$\mathcal{R}(B) = \alpha_{\ell_2} \mathcal{R}_{\ell_2}(B) + \alpha_{tv} \mathcal{R}_{tv}(B) + \alpha_{bn} \mathcal{R}_{bn}(B), \quad (3)$$

where $\mathcal{R}_{\ell_2}(B)$ and $\mathcal{R}_{tv}(B)$ penalize the ℓ_2 norm and the total variation of the recovered watermark data. The detailed forms can be found in Appendix A.1.

Since the attacker \mathcal{A} has no knowledge of the underlying target classes, \mathcal{A} has to conservatively repeat the optimization step in Equation 1 for all possible classes and finally obtain C batches of recovered samples $\{B_c\}_{c=0}^{C-1}$, where B_c indicates the inverted batch of samples towards class c . These recovered samples are expected to contain enough watermark information of f_w .

Note that although the BN layers keep the statistics of the total training set, including samples of all classes, our watermark recovering method is actually using an aggressive class-wise inversion scheme. This is because our goal is to *recover samples close to the watermark data as much as possible here, instead of to generate realistic images*. Since the watermark data may constitute a large contribution to the BN statistics during the training, we propose this scheme to greedily absorb information in BN layers, preventing it from diffusing into other classes. This design is especially effective against watermarks with fixed target classes, compared to the arbitrary-class scheme in conventional inversion methods [61]. We validate this point in Section 6.4.2.

4.2 Watermark Unlearning

Directly unlearning these recovered samples $\{B_c\}_{c=0}^{C-1}$ without other constraints may cause catastrophic forgetting, bringing irreversible significant damage to the normal performance. In the current framework, we consider two data settings for \mathcal{A} 's auxiliary data X_{aux} , in-distribution and transfer. This basic framework is extended to the data-free setting with more improved designs in Section 5.

Assuming the attacker \mathcal{A} has some in-distribution data $X_{aux} = \{(x_i, y_i)\}$ from f_w 's main task, they can leverage these labeled samples to preserve model's performance during unlearning. When \mathcal{A} only has some unlabeled data $\{x_i\}$ from a different distribution, similar goals can be achieved by leveraging model f_w as an oracle to generate pseudo-labels for each sample, yielding the transfer dataset $X_{aux} = \{(x_i, f_w(x_i))\}$.

Starting from f_w , the attacker can finetune the surrogate model f_a under the following optimization objective:

$$\begin{aligned} \min_{\theta_a} \quad & \underbrace{\sum_{(x_i, y_i) \in X_{aux}} \mathcal{L}(f_a(x_i), y_i)}_{\text{preserving original performance}} \\ & + \underbrace{\alpha_{KL} \sum_{B_c \in \{B_c\}_{c=0}^{C-1}} \sum_{\hat{x}_i \in B_c} KL(f_a(\hat{x}_i), y_{soft})}_{\text{removing watermark}}, \end{aligned} \quad (4)$$

where θ_a is the parameters of f_a , $KL(\cdot)$ is the Kullback–Leibler divergence, y_{soft} is a soft unlearning target of length C , where each element equal to $1/C$. Note that we are using a more consistent unlearning target, compared to the hard random labels, which might lead to gradient conflicts and cancellations during fine-tuning.

The first term in Equation 4 stabilizes the general prediction behavior of f_a and thus preserves its performance, while the second term encourages f_a to unlearn the watermark-related information. In this way, we successfully derive a surrogate model f_a to prevent successful watermark verification.

5 IMPROVED DESIGNS FOR DEHYDRA

In this section, we enhance the general framework of DEHYDRA in the following directions.

- (1) **Improved Recovering via Target Class Detection:** First, existing black-box watermarks consist of those with fixed target classes and with non-fixed target classes. For watermarks with fixed target classes (such as *Content* [63] and *EWE* [27]), recovering and unlearning watermark data from other classes bring no benefit to watermark removal, and can even sacrifice the utility of the surrogate model. Therefore, after the original watermark recovering stage, we further distinguish whether the underlying watermark has a fixed target class, and detect its target class if so. We only retain recovered samples belonging to the target class for a fixed-class watermark (§5.1).
- (2) **Improved Unlearning via Recovered Sample Splitting:** Second, the inverted samples derived by Equation 1 may also contain important information of the main task. Indiscriminately unlearning all the inverted samples from the target classes might also hurt the model utility. Therefore, before the original watermark unlearning stage, we split out samples closer to the real watermark data from the recovered ones and improve the finetuning loss accordingly (§5.2).

Figure 2 shows two improved designs and the complete workflow. The two improvements mainly exploit the recovered samples instead of the auxiliary data, enabling our attack to be extended to a data-free setting for watermarks with a fixed target class.

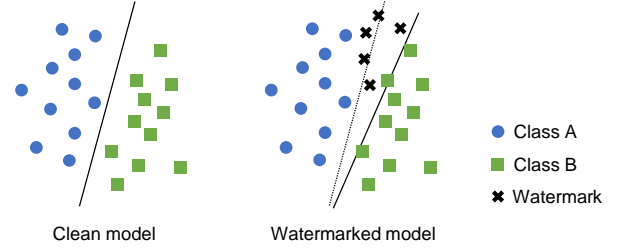


Figure 3: An illustration of the watermark smoothness discrepancy hypothesis. Left: A clean model with its decision boundary. Right: A fixed-class watermark is embedded, with all watermark samples labeled to Class A.

5.1 Target Class Detection

• **Insight: Watermark Smoothness Discrepancy.** By inspecting existing DNN watermarking implementations [2, 27, 33, 42, 63], we find that for a watermark with a fixed target class, the model tends to learn much more diverse data correlated to the target class during training, compared with other labels, due to the additional incorporation of special watermark samples (for the integrity purpose). Moreover, the model also learns samples belonging to the target class more prominently, because of the higher participation ratio of these watermark data (for the robustness purpose).

Based on these observations, we propose the *watermark smoothness discrepancy* hypothesis that, for a DNN model embedded with a fixed-class watermark, the loss landscape (i.e., the structure of loss values around the parameter space) is smoother with respect to samples of the target label, compared to those of other labels. The smoothness property is mathematically defined via the Hessian matrix of loss function [18, 36], which is usually computationally expensive. We instead approximate the smoothness using the performance of model under perturbations. A conceptual illustration is shown in Figure 3, where compared with the clean model, the watermarked model demonstrates higher smoothness on Class A and lower on Class B, due to the additional incorporation of watermark samples and the increased involvement during training.

• **Theoretical Justification.** To provide theoretical evidence, we consider training a watermarked binary classification linear model $f(x) = \text{sign}(\langle w, x \rangle + b)$ on a dataset following mixture Gaussian distributions, together with another set of fixed-class watermark data. The detailed problem definition, theorems and proofs are all deferred to Appendix D due to the page limit.

Our conclusions are (Remark on Theorem D.1): *For a fixed-class watermark, the target class will exhibit a higher smoothness (both input-level and parameter-level) compared to other classes if the watermark samples are (1) diverse internally and different from the normal ones in the target class, and (2) sampled much more frequently. During proof, we also find the additive input-level and parameter-level noise would further enlarge the class-wise discrepancies.*

In existing black-box watermarking implementations, the designs of watermark data for the integrity purpose correspond to the first condition, while the training techniques for the robustness purpose satisfy the second condition. Therefore, we deduce that the watermark smoothness discrepancy hypothesis should also hold for realistic watermarked DNNs.

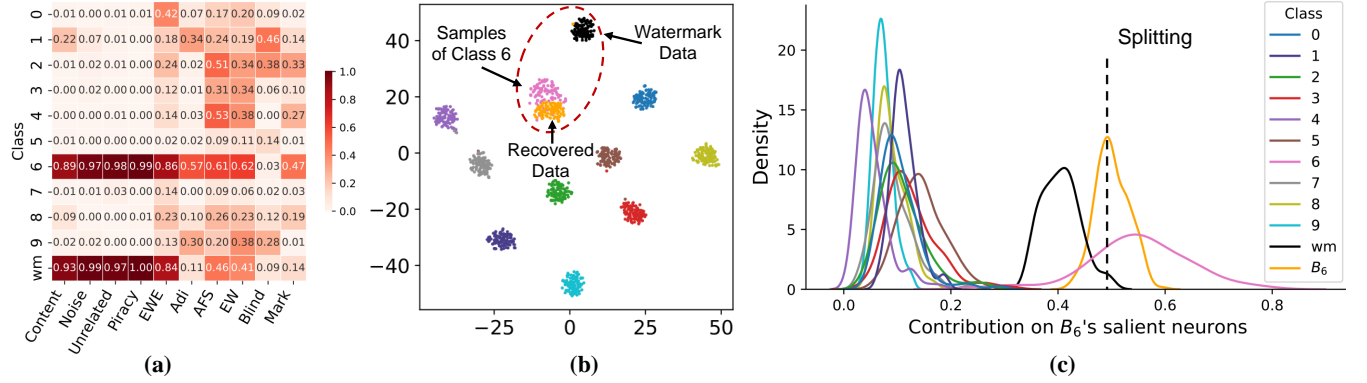


Figure 4: Pilot study results of the improved designs. (a) Class-wise smoothness analysis of the five fixed-class and five non-fixed-class watermarks on CIFAR-10. (b) Activations visualization of a model protected by the *Content* watermark with target class 6. (c) Distribution of sample contributions per class, on the salient neurons of the recovered samples B_6 .

Note that another line of work focuses on class-wise robustness discrepancies [6, 52] and empirically discovers that, for a normal model trained on a balanced dataset, classes with smaller *inter-class* semantic distances are more vulnerable. Our findings are exactly complementary to theirs: the classes with a larger *intra-class* variance and sampling ratio should be more robust.

• **A Pilot Study.** To provide a more comprehensive investigation, we empirically evaluate the smoothness on ten realistic models protected by different black-box watermarks. Following the previous interpretations of smoothness, we define a metric called *SmoothAcc*. Given a set of samples X , we calculate the average prediction accuracy of $(x_i, y_i) \in X$ under perturbations:

$$\text{SmoothAcc}(X) = \frac{1}{|X|} \sum_{(x_i, y_i) \in X} \mathbb{P}(\arg\max_j f_w(x_i + \epsilon_1; \theta_w + \epsilon_2)_j = y_i), \quad (5)$$

where $\epsilon_1 \in \mathcal{N}(0, \sigma_1^2 I_1)$ is the input-level Gaussian noise and $\epsilon_2 \in \mathcal{N}(0, \sigma_2^2 I_2)$ is the parameter-level Gaussian noise.

We construct ten watermarked ResNet-18 [23] models on CIFAR-10 [29] dataset, embedded with five fixed-class watermarks (with the target class set to 6¹) and five non-fixed-class watermarks respectively (more details are described in Section 6.1). For each model, to comprehensively evaluate its smoothness on various data categories, we collect 11 batches of samples, including ten batches of normal samples, each from a class in CIFAR-10, and one batch from real watermark samples, with a batch size of 100. During smoothness analysis, we set the standard deviations of the Gaussian noise $\sigma_1 = 0.5$ and $\sigma_2 = 0.015^2$, and repeatedly sample the Gaussian noise 100 times for each sample to approximate the possibility $\mathbb{P}(\cdot)$ in Equation 5. The results are shown in Figure 4(a).

Evidently, for fixed-class watermarks, the models exhibit much higher smoothness on samples from the target class, including normal samples and watermark samples. This class-wise difference is

less evident for non-fixed-class watermarks, where models achieve moderate smoothness on several classes simultaneously. These results strongly support our hypothesis.

• **Technical Designs.** Based on the observations above, we assume that the watermark smoothness discrepancy phenomenon also exists, when the target model f_w is embedded with a fixed-class watermark, in the class-wise inverted samples $\{B_c\}_{c=0}^{C-1}$. Supposing the fixed target label is y_w , then the class-wise smoothness difference would also be significant in $\{B_c\}_{c=0}^{C-1}$ as the inverted batch B_{y_w} from the target label will cover the space of both normal samples of y_w and watermark samples.

In light of this, we propose to evaluate the smoothness of f_w on $\{B_c\}_{c=0}^{C-1}$ after the watermark recovering, i.e., to calculate $\text{SmoothAcc}(B_c)$ for $c \in \{0, 1, \dots, C-1\}$ respectively. Each $\hat{x}_i \in B_c$ is temporarily labeled as class c during the smoothness evaluation. Next, we sort the batches $\{B_c\}_{c=0}^{C-1}$ according to their $\text{SmoothAcc}(B_c)$ values in the descending order, yielding sorted $\{B_{s_0}, B_{s_1}, \dots, B_{s_{C-1}}\}$. We detect whether f_w is embedded with a fixed-class watermark, and, if so, follow the criterion below to determine the target class: **a)** If the gap between the recovered two batches evaluated with the highest smoothness, i.e., $\text{SmoothAcc}(B_{s_0}) - \text{SmoothAcc}(B_{s_1})$, is larger than a certain threshold T , we regard the underlying watermark as a fixed-class watermark, and take the target label as s_0 . We then only retain recovered samples B_{s_0} and drop other batches. **b)** Otherwise, we regard the underlying watermark as a non-fixed-class watermark, and keep all the recovered batches $\{B_c\}_{c=0}^{C-1}$ since each class c might have several correlated watermark samples. We additionally record labels of the last two batches, s_{C-1} and s_{C-2} , as the least-likely label and second-least-likely label respectively.

In the remainder of this paper, we use $\{B_c\}_*$ to denote $\{B_{s_0}\}$ if a fixed-class watermark is detected, and $\{B_c\}_{c=0}^{C-1}$ otherwise. Consequently, after the watermark recovering and target class detection, we get batch(es) samples $\{B_c\}_*$ close to real watermark data X_w .

5.2 Recovered Samples Splitting

• **Insight: Normal Data Dominance.** We have previously explained the necessity of the class-wise inversion scheme for watermark recovering (§4.1), but we find that this scheme would induce

¹The target classes of the five investigated fixed-class watermarks are set to 6 arbitrarily in this pilot study. We use a randomly chosen label under ten repetitive tests for a systematic study in Section 6.3.

²We at first empirically set the standard deviations of the input-level and parameter-level Gaussian noise. Subsequently, we independently increment each standard deviation until the smoothness discrepancy is evident.

competitive objectives during recovering. Recall the two key components of the recovering objective in Equation 1, i.e. the targeted classification loss $\sum_{\hat{x}_i \in B} \mathcal{L}(f_w(\hat{x}_i), c)$ and the BN regularization term $\sum_{l=1}^L \|\mu_l(B) - \mu_l\|_2^2 + \|\sigma_l^2(B) - \sigma_l^2\|_2^2$. For the recovering process on target class c , the targeted classification loss will lead all the optimized samples B_c consistently classified to c , while the BN regularization term will guide the neuron activations of B_c statistically close to the activations average of the total training set. Considering that the normal samples constitute the main distribution, different from the watermark distribution, the recovered B_c should be closer to class c 's clean data (which belongs to the main distribution) than the watermark samples labeled to c .

Based on the above analysis, we derive another hypothesis named *normal data dominance*, which is intrinsically connected with our class-wise inversion scheme. Specifically, for a watermarked model (with either fixed or non-fixed target classes), the class-wise recovered samples are expected to be closer to normal samples compared with the watermark samples in that class.

• **A Pilot Study.** To gain insights, we analyze the neuron activations of the model protected by *Content* [63] watermark here with fixed target class 6. We first collect 12 batches, including ten batches of normal samples from each class, one batch of watermark samples and one batch of recovered samples B_6 from target class 6, each with a batch size 100. We feed each batch into the watermarked model, extract the activations of the penultimate layer and perform t-SNE dimensionality reduction [55]. As shown in Figure 4(b), the recovered samples (the orange cluster) lie slightly close to the global activations center, and hence closer to the normal samples of class 6 (the pink cluster) than the watermark samples (the black cluster, lying separately from the main distribution). This is because, the watermark task is usually irrelevant to the model's primary task, making it extremely challenging to recover high-quality watermark samples. Although we have employed the class-wise BN regularization term to enhance the coverage of the recovered B_6 over the real watermark samples (§4.1), it would also inevitably guide the recovered samples towards the activation average of the total training set, conforming to our analysis above.

Next, we delve into the composition of the activations of the recovered samples B_6 . Since B_6 are pulled statistically close to the activation center of training data of all classes, it is expected that samples of B_6 will inevitably activate redundant neurons of class c or neurons important to other classes. To prevent the noise from these irrelevant activations, we focus on B_6 's salient neurons. Concretely, we utilize the *importance score* for recognizing the salient neurons [12]. Given a batch of samples X_c of class c , the score extracts the neuron activations at the l -th layer of model f_w , i.e., $f_{w(l)}(X_c)$, and calculates an importance score for each j -th neuron: $\text{Imp}_j(f_{w(l)}(X_c)) = \frac{\mu_j}{\sigma_j}$, where μ_j and σ_j are the mean and the standard deviation of the j -th neuron estimated over X_c . We refer to the neurons with top-5% importance score as salient neurons of X_c , denoted as $SN(X_c)$.

Next, we capture the salient neurons of the recovered samples B_6 and explore the functionality of these consistently activated neurons. For each sample $x_i \in X_c$ among the 12 batches collected above, we define its total activation values on the salient neurons

of B_6 as its contribution:

$$\text{Contribution}(x_i) = \sum_{j \in SN(B_6)} f_{w(l)}(x_i)_j. \quad (6)$$

Then we plot the distribution of sample contributions of each class, distinguished in color. In Figure 4(c), for samples belonging to class 6, normal samples contribute the most, followed by samples in B_6 itself, and finally the watermark samples. This indicates that $SN(B_6)$ are more consistently activated by normal samples of class 6 than the watermark samples, which validates our hypothesis.

More empirical evidence is provided in Figure 8 of Appendix C.1, where, in most cases, $SN(B_6)$ exhibit the highest correlation with the normal samples of class 6. It is worth noting, for non-fixed-class watermarks, $SN(B_6)$ also has a high correlation with the watermark samples besides the normal samples, possibly due to the more complex input-output pairing relations, leaving the watermark data more entangled with the normal data.

• **Technical Designs.** As the class-wise recovered samples B_c might also contain important information of the normal samples in class c , in this section, we exploit the hypothesis of normal data dominance to split B_c into two parts, namely proxy normal data and proxy watermark data. We develop a neuron-level criterion: For each class c , we extract the salient neurons $SN(B_c)$, and take the subset achieving the largest contributions as proxy normal data, leaving the rest as proxy watermark data. Algorithm 1 in Appendix A.2 shows the detailed procedures.

Next, we formally discuss the underlying rationale. An illustrative example is shown in Figure 4(c), where we split the recovered samples B_6 (the orange distribution) according to their contributions on $SN(B_6)$. After splitting, the right part of B_6 are taken as proxy clean data, which are close to the normal samples of class 6 and preserve the normal functionalities specific to class 6. The left part of B_6 are proxy watermark data, either correlating to the real watermark data, or noisily activating other classes' salient neurons (which are inessential for normal performance).

Given the recovered samples $\{B_c\}_*$ obtained via watermark recovering and target class detection, the algorithm will perform class-wise splitting on each batch B_c , returning proxy normal data $\{B_{c,nor}\}_*$ and proxy watermark data $\{B_{c,wmk}\}_*$. Finally, we improve the finetuning loss accordingly.

If the target model is detected with a fixed-class watermark (i.e., the target class s_0), we only need to address samples from class s_0 , denoted as $\{B_{c,wmk}\}_* = \{B_{s_0,wmk}\}$ and $\{B_{c,nor}\}_* = \{B_{s_0,nor}\}$. In this way, the attacker can finetune the surrogate model f_a under the following objective:

$$\begin{aligned} \min_{\theta_a} & \left(\underbrace{\sum_{(x_i, y_i) \in X_{aux}} \mathcal{L}(f_a(x_i), y_i)}_{\text{preserving original performance}} + \sum_{\hat{x}_i \in B_{s_0,nor}} \mathcal{L}(f_a(\hat{x}_i), s_0) \right) \\ & + \underbrace{\sum_{\hat{x}_i \in B_{s_0,wmk}} \mathcal{L}(f_a(\hat{x}_i), \text{rand}_{s_0})}_{\text{removing watermark}}, \end{aligned} \quad (7)$$

where rand_{s_0} denotes a random label in $\{0, 1, \dots, C-1\}$ except s_0 itself. This unlearning target is stronger than minimizing $KL(\cdot, y_{soft})$

in Equation 4, due to the enhanced specificity brought by target class detection, and also better than maximizing $\mathcal{L}(\cdot, s_0)$, due to the easier convergence. Note that in the objective function above, we additionally exploit the proxy normal data $B_{s_0, nor}$ to preserve f_a 's original performance. This design enables our DEHYDRA to be extended to a data-free setting: when the auxiliary data is not available (i.e., the first term enclosed by (\cdot) in the above equation is optional), we can still perform watermark unlearning only using the recovered samples.

For models detected with a non-fixed-class watermark, $\{B_{c, wmk}\}_*$ contains C batches of proxy watermark data. Then \mathcal{A} can finetune f_a under this objective:

$$\min_{\theta_a} \underbrace{\sum_{(x_i, y_i) \in X_{aux}} \mathcal{L}(f_a(x_i), y_i)}_{\text{preserving original performance}} + \underbrace{\sum_{B_{c, wmk} \in \{B_{c, wmk}\}_*} \sum_{\hat{x}_i \in B_{c, wmk}} \mathcal{L}(f_a(\hat{x}_i), y_*)}_{\text{removing watermark}} \quad (8)$$

where y_* denotes the least-likely label s_{C-1} (obtained during the target class detection process) for $B_{c, wmk} \in \{B_{c, wmk}\}_*$ and $c \neq s_{C-1}$, while the second-least-likely label s_{C-2} when $c = s_{C-1}$. This proxy label unlearning target is also stronger than that in Equation 4 due to the enhanced specificity, while maintaining the consistent descending gradients across the labels. Different from the design for fixed-class watermarks, we do not use proxy normal data to further preserve the model performance. This is because the normal data dominance phenomenon is less evident for non-fixed-class watermarks (Figure 8), and thus the salient neurons of the split proxy normal data might still have some intersection with those of the real watermark data. Therefore, leveraging these proxy normal data for preserving utility might hinder the watermark removal.

6 EXPERIMENTS

6.1 Overview of Evaluation

To evaluate the performance of DEHYDRA, we perform a comprehensive study on ten mainstream black-box watermarking schemes, under various benchmark datasets, DNN architectures and data availability settings. Before presenting the detailed evaluation results, we first provide a concise introduction to the experimental setups.

• **Datasets and Victim Models.** Following the settings in [3, 27, 34, 47], we construct victim models on three benchmark datasets, namely, MNIST [32], CIFAR-10 and CIFAR-100 [29]. The respective DNN architectures are LeNet-5 [31], ResNet-18 [23] and ResNet-34. We embed black-box watermarks into the victim models during the training. Then we perform watermark removal attacks, including DEHYDRA and baseline attacks, to evaluate their effectiveness.

• **Target Watermark Schemes.** Our evaluation covers ten mainstream black-box DNN watermarking schemes published at top-tier conferences (some of which are developed by industry leaders such as IBM): *Content*, *Noise*, *Unrelated* [63], *Piracy* [35], *EWE* [27], *Adi* [2], *AFS* [30], *EW* [45], *Blind* [37], *Mark* [20]. These schemes

feature diverse designs in the watermark data and target label settings, and are mostly evaluated in watermark robustness surveys [33, 34, 42]. For more backgrounds, please refer to Appendix B.1.

During implementation, we strictly followed the specifications and hyperparameters in the original papers (with the target class set to 6 for those fixed-class watermarks) to prepare the watermarked models. We also referred to some recent open-source replications [33, 42] to ensure our implementation is faithful. For more implementation details and the performance of the watermarked models, please refer to Appendix B.2 and B.3.

• **Baseline Removal Attacks.** We compare our DEHYDRA with the following three types of removal attacks, comprising six baseline attacks in total.

- **Pruning-based attacks:** (1) **Pruning** [54, 63] directly sets a proportion of DNN parameters with the smallest absolute values to zero. (2) **Fine-pruning** [38] prunes neurons that are infrequently activated by normal data, followed by a finetuning process.
- **Finetuning-based attacks:** (3) **Finetuning** [10, 11] specifically finetunes the target model using a large learning rate, together with a carefully-designed scheduler. (4) **Regularization** [42, 47] finetunes the target model with a large L2 regularization on the parameters. (5) **Distraction** [65] finetunes the target model on in-distribution or transfer data, together with another set of lure data that distracts the model's attention away from the watermark.
- **Unlearning-based attacks:** (6) **Laundering** [3] leverages trigger reverse-engineering methods [56] in backdoor defense literature to recover watermark data, followed by neuron resetting and model retraining, to remove backdoor-based watermarks.

The implementation details of these baseline attacks are clarified in Appendix B.4.

• **Dataset Availability Settings.** We mainly focus on the realistic settings where the attacker has limited data access, including three different scenarios as follows: (1) **In-distribution setting.** The attacker has 1000 correctly-labeled normal samples (2% of the size of the CIFAR-10 training set), following [65]. (2) **Transfer setting.** Here the attacker cannot access the source training data, but is assumed to have 2000 unlabeled samples from another distribution (e.g., CIFAR-100 for target models trained on CIFAR-10, following [21, 65]). (3) **Data-free setting.** In this setting, the attacker cannot access any samples.

• **Implementation of DEHYDRA.** During watermark recovering, we set hyper-parameters $M = 250$, $\alpha_{\ell_2} = 0.01$, $\alpha_{tv} = 0.03$, $\alpha_{bn} = 0.1$. We use $\tanh(\cdot)$ to constrain the batch data within a valid range, and optimize the recovering objective with the Adam optimizer of learning rate 0.1. For target class detection, we set $\sigma_1 = 1.0$, $\sigma_2 = 0.03$ for MNIST and $\sigma_1 = 0.5$, $\sigma_2 = 0.015$ for CIFAR-10 and CIFAR-100 respectively. The detection threshold T is conservatively set to 0.4 for MNIST and CIFAR-10, while 0.3 for CIFAR-100. For splitting the recovered samples, we set l to be the penultimate layer of the target model, and saliency ratio $\beta = 0.95$, split ratio $\gamma = 0.5$ in Alg.1. During watermark unlearning, we set $\alpha_{KL} = 15$ to ensure the unlearning strength for the basic DEHYDRA and the uniform loss weights for the improved DEHYDRA. The model is finally finetuned for 10 epochs using the SGD optimizer with the batch size 128. The learning rate is set to 0.01 when auxiliary data is available

Table 2: Comparison of removal attacks against black-box DNN watermarks under the in-distribution setting on CIFAR-10. x / y denotes the clean accuracy / rescaled watermark accuracy, and values behind \pm report their standard deviations respectively in 5 repetitive tests. The rescaled watermark accuracy values below 50% are bolded and underlined. The left five columns are attack results of the fixed-class watermarks, while the right five are of the non-fixed-class watermarks.

Attacks	Content	Noise	Unrelated	Piracy	EWE	Adi	AFS	EW	Blind	Mark
None	93.8 / 100.0	94.1 / 100.0	93.7 / 100.0	93.3 / 100.0	94.0 / 100.0	93.9 / 100.0	93.0 / 100.0	92.9 / 100.0	91.2 / 100.0	94.0 / 100.0
Fine-pruning	86.4 / <u>44.7</u>	87.4 / 67.0	85.9 / 52.6	69.5 / <u>37.8</u>	88.9 / 70.2	85.4 / 59.8	87.3 / 80.4	87.6 / 83.3	77.0 / 53.8	87.3 / 61.1
Finetuning	85.5 / <u>46.4</u>	85.1 / 84.9	84.6 / <u>47.1</u>	66.1 / <u>40.9</u>	86.9 / 81.2	85.0 / 56.3	90.4 / 99.4	90.1 / 98.8	79.9 / 55.0	85.8 / 57.1
Regularization	60.3 / 100.0	54.7 / 91.5	82.8 / 100.0	76.9 / 100.0	68.2 / <u>47.6</u>	55.4 / 73.8	72.3 / 81.6	69.0 / 82.1	69.8 / 83.6	75.7 / 87.2
Distraction	87.4 / <u>47.5</u>	89.5 / 98.1	88.5 / 51.5	70.6 / <u>40.9</u>	85.8 / 52.4	89.4 / 52.2	92.5 / 99.4	92.5 / 100.0	89.8 / 91.2	86.0 / 51.3
Laundering	91.1 / <u>48.1</u>	90.7 / 77.4	88.8 / 100.0	83.8 / 69.9	91.1 / 70.2	87.0 / 69.7	89.3 / 100.0	89.3 / 98.8	86.3 / 73.1	91.1 / 94.8
DEHYDRABasic	88.8 / <u>44.7</u>	88.1 / <u>5.7</u>	88.1 / <u>44.9</u>	83.3 / <u>44.0</u>	89.0 / <u>48.2</u>	84.1 / 51.6	86.4 / <u>44.9</u>	85.2 / 52.1	84.1 / 58.5	83.0 / 50.7
DEHYDRAImproved	93.1 / <u>45.8</u>	93.1 / <u>5.7</u>	92.7 / <u>44.9</u>	91.4 / <u>41.5</u>	93.3 / <u>47.6</u>	86.0 / <u>49.3</u>	88.1 / <u>44.3</u>	88.2 / <u>46.3</u>	90.1 / <u>46.2</u>	88.5 / <u>47.2</u>

and 0.003 when data-free. Noteworthily, the settings above are empirically chosen and lead to strong attack effectiveness uniformly over almost all the covered watermarking schemes, which supports the watermark-agnostic nature of our DEHYDRA.

• **Evaluation Metrics.** We primarily focus on the utility and the watermark retention of the surrogate models obtained via removal attacks. Specifically, we use *clean accuracy*, i.e., the classification accuracy on the test set, to measure the utility, and *watermark accuracy*, i.e., the ratio of watermark samples correctly classified as target labels, to measure the watermark retention. Further, we follow Lukas et al. [42] to determine the decision threshold θ for each watermark on 20 clean models, and then calculate a metric $S(\cdot; \theta)$ called the *rescaled watermark accuracy*. The former θ is used to distinguish the watermark accuracy of the watermarked model from those independently trained models with high confidence, while the latter $S(\cdot; \theta)$ enables a fair robustness comparison over different watermarking schemes. Specially, the rescaling function is defined as $S(x; \theta) = \max(0, \frac{1-\theta'}{1-\theta}x + \frac{\theta'-\theta}{1-\theta})$, linearly rescaling the watermark accuracy x based on the decision threshold θ . By manually picking a rescaled decision threshold $\theta' = 0.5$, a watermark is said to be removed if the rescaled watermark accuracy is lower than 50%. Table 12 in Appendix B.3 summarizes our estimated decision threshold for the ten black-box watermarking schemes.

We define a removal attack as successful (i.e., cracks the target watermark scheme) if the rescaled watermark accuracy is below 50% and the surrogate model maintains at least 90% of the original utility. Note that this criterion is slightly different from [42] where a maximum 5% utility loss is used. This is mainly because we consider the more realistic data-limited settings (e.g., 2% of the source training set in most experiments), while [42] assume the attacker has access to over 30% of the training set. Nevertheless, we also evaluate the attack performance when more data is available in Section 7.

6.2 Attack Performance

The comparison with baseline removal attacks is organized according to the data settings where the baseline attacks are applicable. Due to the page limit, we mainly present and analyze the attack results on CIFAR-10, while similar results are observed on MNIST and CIFAR-100. Appendix C.2 and C.3 present the omitted results. (1) **In-distribution Setting.** Table 2 presents the attack results by performing different removal attacks against the ten investigated

watermarks under the in-distribution setting. Both the basic DEHYDRA and improved DEHYDRA significantly lower the rescaled watermark accuracy of surrogate models, among which the improved DEHYDRA achieves comprehensive watermark removal against the ten investigated black-box watermarks, with at least 90% utility preserved (six in ten even with a clean accuracy degradation in 2%).

The baseline removal attacks are only effective against a subset of watermark schemes. For instance, *Laundering* is more effective against watermarks similar to backdoors, such as *Content*, *Piracy* and *EWE*, but less effective against other watermarks, because of its dependence on the trigger reverse-engineering method to detect and remove the underlying watermarks. *Fine-pruning* and *Finetuning* are less effective against *EWE*, *Adi* and *AFS*, possibly due to the more entangled learned representations of the watermark data. *Distraction* can also only compromise about half of these watermarks' robustness, while being less powerful against watermarks of other forms. In contrast, our basic and improved DEHYDRA are both effective regardless of the target watermark schemes, due to the capture of the black-box watermarks' commonality. Note that our improved DEHYDRA can effectively remove the *Piracy* watermark, despite the fact that the number of its watermark data (1% of the training data, i.e., 500, following the original paper) is much larger than the number of unlearned samples in our improved DEHYDRA ($\gamma \times M = 125$). This indicates that the removal effect of our approach is achieved not just through a sample-level unlearning, but also through a deeper neuron-level fixing.

Additionally, these baseline attacks might impact the surrogate model's utility seriously. For instance, the finetuning-based attacks, such as *Finetuning*, *Regularization* and *Distraction*, incur a large clean accuracy degradation against *Noise* and *Piracy* watermarks, possibly due to their overly offensive finetuning strategies. However, our DEHYDRA attacks could still preserve the surrogate model's performance besides removing the watermark, because of their more specific training (unlearning) objectives.

(2) **Transfer Setting.** As we can see from Table 3, the clean accuracy of the surrogate models are generally lower than the in-distribution setting, since the proxy labels of these transfer data may contain noise. *Regularization* is generally ineffective, and *Distraction* is only effective against a subset of fixed-class watermarks.

On the contrary, both the basic and improved DEHYDRA significantly lower the rescaled watermark accuracy while preserving the clean accuracy at an acceptable level. The improved DEHYDRA

Table 3: Comparison of removal attacks against black-box watermarks under the transfer setting on CIFAR-10.

Attacks	Content	Noise	Unrelated	Piracy	EWE	Adi	AFS	EW	Blind	Mark
None	93.8 / 100.0	94.1 / 100.0	93.7 / 100.0	93.3 / 100.0	94.0 / 100.0	93.9 / 100.0	93.0 / 100.0	92.9 / 100.0	91.2 / 100.0	94.0 / 100.0
Regularization	49.3 / 99.4	60.3 / 73.6	53.1 / 100.0	60.5 / 100.0	53.8 / 47.6	41.0 / 69.7	43.4 / 71.5	55.5 / 76.9	44.6 / 69.6	53.4 / 74.5
Distraction	89.0 / 50.9	89.6 / 100.0	88.7 / 49.8	84.2 / 49.0	88.8 / 63.4	87.7 / 60.3	92.0 / 100.0	92.3 / 100.0	88.8 / 91.2	89.1 / 66.9
DEHYDRA _{Basic}	85.5 / 49.8	85.3 / 6.7	80.0 / 46.0	82.4 / 45.2	88.2 / 50.8	80.8 / 54.5	84.6 / 48.1	82.1 / 58.5	80.2 / 62.0	84.0 / 55.9
DEHYDRA _{Improved}	92.5 / 47.0	92.3 / 5.7	92.7 / 48.2	90.1 / 47.1	92.6 / 49.7	80.4 / 59.8	86.2 / 47.4	83.5 / 49.8	88.8 / 46.8	87.0 / 49.5

Table 4: Comparison of removal attacks against black-box watermarks under the data-free setting on CIFAR-10.

Attacks	Content	Noise	Unrelated	Piracy	EWE
None	93.8 / 100.0	94.1 / 100.0	93.7 / 100.0	93.3 / 100.0	94.0 / 100.0
Pruning	81.5 / 78.2	90.7 / 42.5	78.7 / 45.4	41.9 / 37.8	92.2 / 100.0
DEHYDRA (Improved)	89.8 / 45.8	88.6 / 5.7	88.8 / 46.0	90.0 / 39.6	90.7 / 48.2

still successfully cracks nine of the ten watermarks under this setting. Note that our DEHYDRA cannot completely remove the *Adi* watermark (despite the rescaled watermark accuracy being lower than 60%) under the transfer setting. This is reasonable because *Adi* uses abstract out-of-distribution images as watermark data, whose salient neurons are more likely to be activated by the transfer dataset, compared to other common watermarks.

(3) Data-free Setting. Almost no baseline attacks except for the *Pruning* attack are applicable to the data-free scenario. Therefore, as also explained in Section 5.2, we mainly compare the improved DEHYDRA, when attacking five fixed-class watermarks here, with *Pruning*. The attack results are shown in Table 4.

As is shown, our improved DEHYDRA achieves surprisingly good performance in the data-free setting, significantly surpassing the baseline *Pruning*. All the five fixed-class watermarks are completely removed from the protected model, with a clean accuracy degradation of no more than 5.5%. This validates the correctness of our splitting algorithm on the recovered samples, since the surrogate modes are finetuned solely on the recovered samples $\{B_c\}_* = \{B_{s_0}\}$, with split proxy normal data $\{B_{s_0,nor}\}$ to maintain the utility, and proxy watermark data $\{B_{s_0,wmk}\}$ for removal, as in Equation 7.

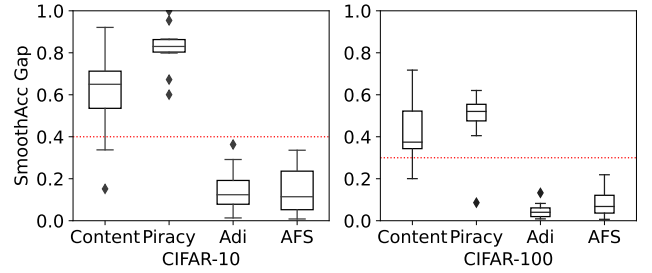
6.3 Validating Detection of Target Classes

In this part, we mainly study the risks of target class detection errors. We temporarily ignore the case when fixed-class watermarks are correctly identified but given the incorrect target label, since this is extremely rare, as validated by later quantitative experiments. Here we mainly focus on two types of errors, false positives (non-fixed-class watermarks identified as fixed-class ones) and false negatives (fixed-class watermarks identified as non-fixed-class ones).

Intuitively, we have a lower tolerance for false positives compared to false negatives. As shown in Table 5, for fixed-class watermarks (*Content* and *Piracy*) wrongly identified as non-fixed-class ones, the clean accuracy will degrade a little (e.g., within 3% on CIFAR-10) and the watermark is still effectively removed. However, for non-fixed-class watermarks (*Adi* and *AFS*) falsely identified as fixed-class ones, the removal effectiveness will be significantly hampered (e.g. the rescaled watermark accuracy remains over 90%). Therefore, we should set the detection threshold conservatively to prevent false positive cases.

Table 5: Attack results of improved DEHYDRA on CIFAR-10 and CIFAR-100, with the target class detection results deliberately *resp.* set to fixed-class and non-fixed-class ones.

Dataset	Detection	Content	Piracy	Adi	AFS
CIFAR-10	fixed	93.1 / 45.8	91.4 / 41.5	92.0 / 96.5	88.8 / 93.7
	non-fixed	90.7 / 44.2	88.9 / 40.3	86.0 / 49.3	88.1 / 44.3
CIFAR-100	fixed	75.1 / 48.7	71.0 / 41.8	64.0 / 71.5	67.1 / 100.0
	non-fixed	67.7 / 48.7	65.8 / 40.0	64.8 / 60.3	65.4 / 54.8

**Figure 5: SmoothAcc gap distribution under ten random tests. The box plot shows min/max and quartiles, as well as the estimated outliers.**

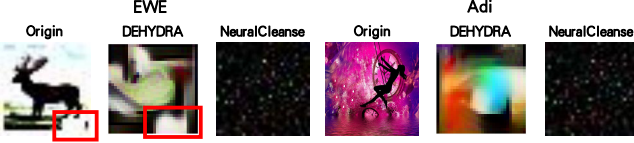
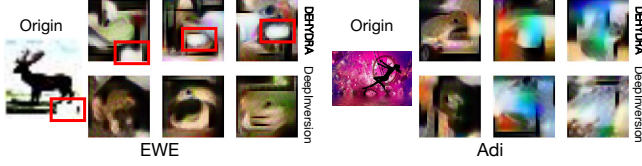
Next, we further investigate the detection accuracy of our target class detection algorithm. For each watermark scheme on each dataset, we independently train ten watermarked models, each using the randomly re-generated watermark data with different target classes. For each target watermarked model, we perform watermark recovering and then target class detection. We calculate the SmoothAcc gap $\text{SmoothAcc}(B_{y_w}) - \max_{i \neq y_w} \text{SmoothAcc}(B_i)$ for a fixed-class watermark with target class y_w , and $\text{SmoothAcc}(B_{s_1}) - \text{SmoothAcc}(B_{s_2})$ for a non-fixed-class watermark. As shown in Figure 5, the SmoothAcc gaps for fixed-class watermarks are always positive in the ten tests, which means we can always determine a correct target class for those successfully identified fixed-class watermarks. Besides, the SmoothAcc gaps of fixed-class and non-fixed-class watermarks are large enough to find a threshold to determine them, which validates the correctness of our target class detection algorithm.

6.4 Ablation Studies

6.4.1 Effectiveness of the Two Improvements. We further validate the effectiveness of the two proposed improvements. As shown in Table 6, target class detection is especially useful to fixed-class watermarks. It is also useful to non-fixed-class watermarks, due to the easier unlearning target of proxy hard labels. The splitting algorithm is generally beneficial to both fixed-class and non-fixed-class watermarks to preserve normal performance, due to the enhanced

Table 6: Ablation study on the two improvements, target class detection and recovered samples splitting, under the in-distribution setting on CIFAR-10.

Attacks	Unrelated	EWE	EW	Mark
DEHYDRA _{Basic}	88.1 / 44.9	89.0 / 48.2	85.2 / 52.1	83.0 / 50.7
DEHYDRA + Detection	91.9 / 44.9	92.4 / 47.6	86.8 / 46.9	87.8 / 48.4
DEHYDRA + Splitting	88.3 / 44.9	89.6 / 47.6	87.2 / 51.5	85.7 / 48.9
DEHYDRA _{Improved}	92.7 / 44.9	93.3 / 47.6	88.2 / 46.3	88.5 / 47.2

**Figure 6: Visualization of the recovered watermark images/triggers by DEHYDRA/NeuralCleanse.****Figure 7: Visualization of the recovered samples by DEHYDRA and DeepInversion.****Table 7: Attack results using different recovering methods on CIFAR-10.**

Inversion Method	Content	EWE	Adi
DEHYDRA	92.6 / 44.7	92.5 / 48.2	85.1 / 48.7
NeuralCleanse	92.5 / 47.5	92.2 / 98.4	84.4 / 93.0
DeepInversion	91.3 / 68.7	91.2 / 99.0	84.3 / 49.3

removal specificity. Armed with the two designs, the improved DEHYDRA achieves the highest clean accuracy, with a 4% increase in average compared with the basic version.

6.4.2 Comparison with Other Recovering Methods. We further prove the effectiveness of the watermark recovering algorithm of DEHYDRA, by comparing with representative algorithms from (1) **Trigger Reverse-Engineering** (e.g., an adaptive version of *NeuralCleanse* [56]), which is originally designed to detect and erase backdoors in DNNs, and (2) **General Model Inversion** (e.g., *DeepInversion* [61]), which is designed to synthesize realistic training samples of the target models for privacy disclosure or knowledge distillation. Our experiments below focus on the in-distribution data setting on CIFAR-10, and perform watermark removal attacks against three black-box watermarks, *Content*, *EWE* and *Adi*. More implementation details are deferred to Appendix B.5.

Visual Effects. Figure 6 visualizes the recovered watermark data by DEHYDRA and the trigger patterns by *NeuralCleanse*, while Figure 7 compares DEHYDRA with *DeepInversion*. As is shown, the recovered triggers by *NeuralCleanse* are visually meaningless and rather similar to each other. This implies that the assumed shortcut in the model’s decision space is non-existent in these watermarked

Table 8: Effectiveness of DEHYDRA on watermarked models after differential privacy training.

Eps	Content		Adi	
	Before	After	Before	After
$\epsilon = 1$	57.3 / 100.0	51.5 / 100.0	53.0 / 100.0	45.0 / 68.5
$\epsilon = 2$	62.0 / 100.0	55.9 / 100.0	60.0 / 99.4	53.9 / 62.1
$\epsilon = 6$	79.8 / 100.0	63.4 / 83.3	77.8 / 99.4	62.8 / 57.4

Table 9: The effectiveness of deceiving the detection algorithm against our improved DEHYDRA.

N_{trap}	Adi			EW		
	Before	After	Gap	Before	After	Gap
100	87.6 / 100.0	84.4 / 60.3	0.1233	90.2 / 100.0	87.4 / 51.5	0.0766
150	90.3 / 100.0	85.1 / 60.9	0.1291	92.7 / 100.0	89.9 / 92.5	0.7825
200	89.4 / 100.0	89.5 / 95.9	0.5406	90.3 / 100.0	90.0 / 98.3	0.6226

Table 10: The effectiveness of bypassing the splitting algorithm against our improved DEHYDRA.

WM size	Content		EWE	
	Before	After	Before	After
1000	91.9 / 100.0	89.0 / 52.0	90.3 / 99.5	89.2 / 62.3
2000	91.8 / 100.0	90.3 / 53.1	90.0 / 100.0	89.2 / 74.9
5000	88.5 / 100.0	87.9 / 58.1	89.4 / 100.0	88.6 / 79.6

models, and the recovered triggers are only the universal adversarial perturbations of the CIFAR-10 dataset itself. Also, *DeepInversion* tends to produce more realistic images belonging to the target classes (e.g., “frogs” for *EWE*), but less informative of the watermark data. In comparison, most of the recovered watermark data holds the essential features of the real watermark patterns, such as the white patches for *EWE* and colorful fringing for *Adi*. These visual differences demonstrate the effectiveness of our aggressive inversion scheme.

Attack Results. More quantitatively, we investigate the removal effectiveness using the three recovering methods, with attack results shown in Table 7. As we can see, *NeuralCleanse* is ineffective against *EWE* and *Adi*, due to the inexistence of backdoor shortcuts. For *DeepInversion*, unlearning the inverted samples from class 6 fails to remove *EWE*, because the BN statistics of the watermark data might diffuse into other classes during its random-class inversion scheme. The recovered samples are therefore more representative of the normal data, unlearning which incurs a larger utility loss. Only DEHYDRA is effective for removing all three watermarks from the protected models, conforming to the previous visual observations.

6.5 Potential Defenses

The experiments above validate that DEHYDRA poses real threats to existing black-box watermarks. In this section, we further consider potential defenses against DEHYDRA. Specifically, we design defenses to resist or bypass the following three key components of DEHYDRA: the inversion-based framework, target class detection, and recovered samples splitting.

Differential Privacy. Differential privacy (DP) is a general approach for privacy protection [1]. Here, we investigate the effectiveness of differential privacy to reduce the watermark information

Table 11: Comparison to extraction attacks under different data ratios w.r.t. the training data on CIFAR-10.

Data Ratio	Noise			EWE			AFS		
	DEHYDRA	Extraction	Combined	DEHYDRA	Extraction	Combined	DEHYDRA	Extraction	Combined
2%	93.1 / 5.7	46.5 / 24.6	45.5 / 14.2	93.3 / 47.6	48.8 / 56.5	46.6 / 50.3	88.1 / 44.3	48.0 / 44.9	45.8 / 45.5
5%	93.3 / 10.4	62.1 / 88.7	56.6 / 5.7	93.1 / 48.7	58.7 / 55.0	56.0 / 49.2	90.4 / 46.2	63.5 / 46.2	53.5 / 48.1
10%	93.4 / 8.5	72.4 / 18.0	65.6 / 5.7	93.3 / 49.7	70.1 / 48.2	66.3 / 49.2	91.2 / 46.8	71.8 / 46.8	67.4 / 48.7
33.33%	94.0 / 8.5	86.3 / 81.1	78.5 / 33.1	93.2 / 49.2	81.8 / 60.2	78.9 / 47.6	92.8 / 50.6	86.7 / 45.5	78.5 / 50.0

inverted by DEHYDRA. Roughly speaking, we use DP-SGD with the privacy budgets (ϵ, δ) during the watermark embedding process: δ is set to be a fixed value of 1×10^{-5} and the ϵ is varied to evaluate the attack effect under different privacy budgets, a smaller ϵ means adding larger noise to the gradients. As we can see in Table 8, with the decrease of ϵ , DEHYDRA will achieve a worse watermark removal. This is because the noise added to the gradients during the training phase makes the learned features of watermark information more scattered, thus compromising the effectiveness of the attack. However, this is brought with a significant accuracy drop of 40% on average.

Deceiving the Detection Algorithm. Here we try to cause the severer false positive errors (making non-fixed-class watermarks detected as fixed-class ones) in Section 6.3. Concretely, model owners can deliberately change the watermark smoothness by adding N_{trap} watermark samples with the specified identical labels (in addition to the original 100 watermark samples with non-fixed target classes) during training. In Table 9, when $N_{trap} = 200$, the SmoothAcc Gap will be larger than the detection threshold, and the watermark will be detected as fixed-class ones. However, forcing models to memorize more samples will cause an accuracy drop about 3% on average.

Bypassing the Splitting Algorithm. The model owner may attempt to bypass the recovered samples splitting algorithm, by deliberately adding more watermark samples into the target classes and disrupting the dominance of normal data. This may decrease the attack effectiveness of the improved DEHYDRA in data-limited settings, as we can see in Table 10, the larger trigger set size, the worse the removal effect will be. However, the attacker may skip the splitting procedure and restore the attack effectiveness with in-distribution/transfer data.

7 DISCUSSION

Black-box DNN Watermarks and Backdoors. In previous literature, black-box watermarks are sometimes also called backdoor-based [3, 39, 47] or trigger set-based watermarks [34]. However, black-box watermarks have substantial differences from DNN backdoors in terms of the watermark data patterns and target class settings (§2.1), especially for the non-fixed-class watermarks.

Therefore, attacking black-box watermarks has its unique challenges, especially for an attacker who usually has no knowledge of the underlying watermarking schemes. This also explains why existing removal attacks extended from backdoor defenses, such as *Laundrying* (§6.2) and *NeuralCleanse* (§6.4.2), could not effectively remove most of the watermarks.

Attack Performance when More Data is Available. Our work primarily focuses on the data-limited settings, where only 2% of the training data is used in most of the experiments (§6.2). This

limited data availability may have a slight impact on the utility of the surrogate model for several non-fixed-class watermarks. In this part, we evaluate the performance of the improved DEHYDRA when more data is available. As shown in Table 11, for the non-fixed-class watermark AFS, the utility loss is mitigated when more data is available. The accuracy drop is controlled within 1% for the three investigated watermarks, when 33% of the training data is available.

Comparison to Model Extraction Attacks. Model extraction attacks are effective for watermark removal but typically require a large query dataset. Here we compare the improved DEHYDRA with the classical *Extraction* attack [27] and the recently proposed *Combined* attack (which combines transfer learning with label smoothing) that has been reported as effective against all watermarks using one-third of the training data [42]. As shown in Table 11, the surrogate models obtained through extraction attacks demonstrate significantly dependent performance on the available dataset size, resulting in relatively low accuracy when the data ratio is 2%, 5%, and 10%. In contrast, our DEHYDRA maintains high accuracy while effectively removing the watermark across all data ratios.

Comparison to Data Extraction Attacks. Data extraction attacks aim at reconstructing training data from a well-trained DNN model. Since DNN watermarks are typically designed to trace ownership after model leakage (§2.1), it is essential to thoroughly assess the privacy risk of the watermark data when the model parameters are completely exposed. In fact, previous *Content* [63] and *EWE* [27] watermarks did investigate the watermark privacy risk, while using naive extraction methods, leading to the conclusion that their watermarks are "secure". On the contrary, DEHYDRA specially targets the overfitted watermark data, and hence adopts an aggressive class-wise inversion method to better exploit model internals. Ablation studies in Section 6.4.2 also demonstrate that our DEHYDRA captures more meaningful visual features of the watermark samples (Figure 6, 7) and achieves significantly better removal results (Table 7) compared to traditional data extraction methods.

Attack Efficiency. For a target model with a large number of class labels, the main computation cost of DEHYDRA comes from the class-wise watermark recovering process, which is proportional to the number of classes. Similar to the solutions in *NeuralCleanse* [56], we find parallel computation on multiple GPUs and early-stopping mechanisms also useful here. Moreover, we propose a multi-class inversion mechanism while optimizing a mixed data batch, where the BN term is calculated for samples of each class separately.

Limitations and Future Works. Our methodological design and evaluations primarily focus on the image classification task, which is consistent with the mainstream DNN watermarking research [33, 34, 42]. We note that there do exist some recent works extending black-box DNN watermarks to other domains, including

self-supervised learning for encoders, graph-based tasks, and text-based tasks. For these new applications, we posit that the watermark information is still deeply embedded within the internals of the watermarked models, suggesting that DEHYDRA continues to be a promising and effective attack. The primary challenge to adapt the DEHYDRA attack lies in developing a suitable recovering method (for instance, prompting-based methods [8] could be utilized to recover the watermark information stored in a large language model). Nonetheless, we leave it as future work to explore the further effectiveness of DEHYDRA on other tasks.

8 CONCLUSION

In this paper, we identify a shared trait of existing black-box DNN watermarks: they exploit the over-parameterization of DNNs to especially memorize the watermark data correlated to the target classes. Based on this, we propose a novel watermark-agnostic removal framework, namely DEHYDRA, which exploits the model internals to recover and unlearn the underlying watermarks. With in-depth analysis, theoretical justifications and pilot studies, we further improve DEHYDRA with target class detection and recovered sample splitting algorithms, which help preserve the model utility and significantly reduce the dataset dependence. Extensive evaluations on various settings against ten mainstream black-box watermarks demonstrate that DEHYDRA is generally effective and has minimal utility impact, with much more relaxed or even no assumptions on the dataset availability.

ACKNOWLEDGMENTS

We would like to thank the anonymous reviewers for their insightful comments that helped improve the quality of the paper. This work was supported in part by the National Key Research and Development Program (2021YFB3101200), National Natural Science Foundation of China (62472096, 62172104, 62172105, 62102093, 62102091, 62302101, 62202106). Min Yang is a faculty of Shanghai Institute of Intelligent Electronics & Systems and Engineering Research Center of Cyber Security Auditing and Monitoring, Ministry of Education, China. Mi Zhang and Min Yang are the corresponding authors.

REFERENCES

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*. 308–318.
- [2] Yossi Adi, Carsten Baum, Moustapha Cisse, Benny Pinkas, and Joseph Keshet. 2018. Turning your weakness into a strength: Watermarking deep neural networks by backdoor. In *27th {USENIX} Security Symposium ({USENIX} Security 18)*. 1615–1631.
- [3] William Aiken, Hyoungshick Kim, and Simon S. Woo. 2020. Neural Network Laundering: Removing Black-Box Backdoor Watermarks from Deep Neural Networks. *Comput. Secur.* 106 (2020), 102277.
- [4] Marco Allodi, Alberto Broggi, Domenico Giaquinto, Marco Patander, and Antonio Prioletti. 2016. Machine learning in tracking associations with stereo vision and lidar observations for an autonomous vehicle. In *2016 IEEE intelligent vehicles symposium (IV)*. IEEE, 648–653.
- [5] Mrinal R Bachute and Javed M Subhedar. 2021. Autonomous driving architectures: insights of machine learning and deep learning algorithms. *Machine Learning with Applications* 6 (2021), 100164.
- [6] Philipp Benz, Chaoning Zhang, Adil Karjauv, and In So Kweon. 2021. Robustness may be at odds with fairness: An empirical study on class-wise accuracy. In *NeurIPS 2020 Workshop on Pre-registration in Machine Learning*. PMLR, 325–342.
- [7] Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. 2019. The secret sharer: Evaluating and testing unintended memorization in neural networks. In *28th USENIX security symposium (USENIX security 19)*. 267–284.
- [8] Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, et al. 2021. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*. 2633–2650.
- [9] Xuxi Chen, Tianlong Chen, Zhenyu Zhang, and Zhangyang Wang. 2021. You are caught stealing my winning lottery ticket! Making a lottery ticket claim its ownership. *Advances in Neural Information Processing Systems* 34 (2021), 1780–1791.
- [10] Xinyun Chen, Wenxiao Wang, Chris Bender, Yiming Ding, Ruoxi Jia, Bo Li, and Dawn Song. 2021. Refit: a unified watermark removal framework for deep learning systems with limited data. In *Proceedings of the 2021 ACM Asia Conference on Computer and Communications Security*. 321–335.
- [11] Xinyun Chen, Wenxiao Wang, Yiming Ding, Chris Bender, Ruoxi Jia, Bo Li, and Dawn Song. 2019. Leveraging unlabeled data for watermark removal of deep neural networks. In *ICML workshop on Security and Privacy of Machine Learning*. 1–6.
- [12] Zhenzhu Chen, Shang Wang, Anmin Fu, Yansong Gao, Shui Yu, and Robert H. Deng. 2022. LinkBreaker: Breaking the Backdoor-Trigger Link in DNNs via Neurons Consistency Check. *IEEE Transactions on Information Forensics and Security* 17 (2022), 2000–2014.
- [13] Gregory Cohen, Saeed Afshar, Jonathan Tapson, and Andre Van Schaik. 2017. EMNIST: Extending MNIST to handwritten letters. In *2017 international joint conference on neural networks (IJCNN)*. IEEE, 2921–2926.
- [14] The Turing Way Community. 2021. *Licensing Machine Learning models*. The Turing Way. <https://book.the-turing-way.org/reproducible-research/licensing/licensing-ml.html> Accessed on July 21, 2024.
- [15] Bitu Darvish Rouhani, Huili Chen, and Farinaz Koushanfar. 2019. Deepsigns: An end-to-end watermarking framework for ownership protection of deep neural networks. In *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*. 485–497.
- [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [17] Hugging Face. 2024. *Hugging Face Models*. Retrieved April 10, 2024 from <https://huggingface.co/models>
- [18] Zhijian Ge, Fanhua Shang, Hongying Liu, Yuan Yuan Liu, and Xiaosen Wang. 2023. Boosting Adversarial Transferability by Achieving Flat Local Maxima. *arXiv preprint arXiv:2306.05225* (2023).
- [19] Tristan Greene. 2018. *IBM came up with a watermark for neural networks*. Retrieved April 10, 2024 from <https://thenextweb.com/news/ibm-came-up-with-a-watermark-for-neural-networks>
- [20] Jia Guo and Miodrag Potkonjak. 2018. Watermarking deep neural networks for embedded systems. In *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*. IEEE, 1–8.
- [21] Shangwei Guo, Tianwei Zhang, Han Qiu, Yi Zeng, Tao Xiang, and Yang Liu. 2020. Fine-tuning Is Not Enough: A Simple yet Effective Watermark Removal Attack for DNN Models. In *International Joint Conference on Artificial Intelligence*.
- [22] Abdalraouf Hassan and Ausif Mahmood. 2018. Convolutional recurrent deep learning model for sentence classification. *Ieee Access* 6 (2018), 13949–13957.
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [24] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4700–4708.
- [25] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and Aleksander Madry. 2019. Adversarial examples are not bugs, they are features. *Advances in neural information processing systems* 32 (2019).
- [26] Sergey Ioffe and Christian Szegedy. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*. pmlr, 448–456.
- [27] Hengrui Jia, Christopher A Choquette-Choo, Varun Chandrasekaran, and Nicolas Papernot. 2021. Entangled Watermarks as a Defense against Model Extraction. In *USENIX Security Symposium*. 1937–1954.
- [28] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. 2017. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences* 114, 13 (2017), 3521–3526.
- [29] Alex Krizhevsky, Geoffrey Hinton, et al. 2009. Learning multiple layers of features from tiny images. (2009).
- [30] Erwan Le Merrer, Patrick Perez, and Gilles Trédan. 2020. Adversarial frontier stitching for remote neural network watermarking. *Neural Computing and Applications* 32 (2020), 9233–9244.
- [31] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324.

- [32] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. 1998. The MNIST Database of Handwritten Digits. <http://yann.lecun.com/exdb/mnist/> (1998).
- [33] Isabell Lederer, Rudolf Mayer, and Andreas Rauber. 2023. Identifying Appropriate Intellectual Property Protection Mechanisms for Machine Learning Models: A Systematization of Watermarking, Fingerprinting, Model Access, and Attacks. *arXiv preprint arXiv:2304.11285* (2023).
- [34] Suyoung Lee, Wonho Song, Suman Jana, Meeyoung Cha, and Soeul Son. 2022. Evaluating the robustness of trigger set-based watermarks embedded in deep neural networks. *IEEE Transactions on Dependable and Secure Computing* (2022).
- [35] Huiying Li, Emily Wenger, Shawn Shan, Ben Y Zhao, and Haitao Zheng. 2019. Piracy resistant watermarks for deep neural networks. *arXiv preprint arXiv:1910.01226* (2019).
- [36] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. 2018. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems* 31 (2018).
- [37] Zheng Li, Chengyu Hu, Yang Zhang, and Shanqing Guo. 2019. How to prove your model belongs to you: a blind-watermark based framework to protect intellectual property of DNN. *Proceedings of the 35th Annual Computer Security Applications Conference* (2019). <https://api.semanticscholar.org/CorpusID:207847538>
- [38] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. 2018. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *Research in Attacks, Intrusions, and Defenses: 21st International Symposium, RAID 2018, Heraklion, Crete, Greece, September 10–12, 2018, Proceedings 21*. Springer, 273–294.
- [39] Xuankai Liu, Fengting Li, Bihan Wen, and Qi Li. 2021. Removing backdoor-based watermarks in neural networks with limited data. In *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 10149–10156.
- [40] Yifan Lu, Wenxuan Li, Mi Zhang, Xudong Pan, and Min Yang. 2024. Neural Dehydration: Effective Erasure of Black-box Watermarks from DNNs with Limited Data. *arXiv preprint arXiv:2309.03466* (2024).
- [41] Nils Lukas. 2022. *Watermark-Robustness-Toolbox*. Retrieved April 10, 2024 from <https://github.com/dnn-security/Watermark-Robustness-Toolbox>
- [42] Nils Lukas, Edward Jiang, Xinda Li, and Florian Kerschbaum. 2022. Sok: How robust is image classification deep neural network watermarking?. In *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE, 787–804.
- [43] Kinsong Ma, Zekai Wang, and Weiwei Liu. 2022. On the tradeoff between robustness and fairness. *Advances in Neural Information Processing Systems* 35 (2022), 26230–26241.
- [44] Alexander Mordvintsev, Christopher Olah, and Mike Tyka. 2015. Inceptionism: Going deeper into neural networks. (2015).
- [45] Ryota Namba and Jun Sakuma. 2019. Robust watermarking of neural network with exponential weighting. In *Proceedings of the 2019 ACM Asia Conference on Computer and Communications Security*. 228–240.
- [46] Omkar M. Parkhi, Andrea Vedaldi, and Andrew Zisserman. 2015. Deep Face Recognition. In *British Machine Vision Conference*. <https://api.semanticscholar.org/CorpusID:4637184>
- [47] Masoumeh Shafieinejad, Jiaqi Wang, Nils Lukas, and Florian Kerschbaum. 2019. On the Robustness of Backdoor-based Watermarking in Deep Neural Networks. *Proceedings of the 2021 ACM Workshop on Information Hiding and Multimedia Security* (2019).
- [48] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).
- [49] Shichang Sun, Haoqi Wang, Mingfu Xue, Yushu Zhang, Jian Wang, and Weiqiang Liu. 2021. Detect and Remove Watermark in Deep Neural Networks via Generative Adversarial Networks. In *Information Security - 24th International Conference, ISC 2021, Virtual Event, November 10–12, 2021, Proceedings (Lecture Notes in Computer Science, Vol. 13118)*, Joseph K. Liu, Sokratis K. Katsikas, Weizhi Meng, Willy Susilo, and Rolly Intan (Eds.). Springer, 341–357.
- [50] Zhichuang Sun, Ruimin Sun, Long Lu, and Alan Mislove. 2021. Mind Your Weight(s): A Large-scale Study on Insufficient Machine Learning Model Protection in Mobile Apps. In *30th USENIX Security Symposium, USENIX Security 2021, August 11–13, 2021*, Michael D. Bailey and Rachel Greenstadt (Eds.). USENIX Association, 1955–1972. <https://www.usenix.org/conference/usenixsecurity21/presentation/sun-zhichuang>
- [51] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2818–2826.
- [52] Qi Tian, Kun Kuang, Kelu Jiang, Fei Wu, and Yisen Wang. 2021. Analysis and applications of class-wise robustness in adversarial training. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 1561–1570.
- [53] Florian Tramèr, Fan Zhang, Ari Juels, Michael K Reiter, and Thomas Ristenpart. 2016. Stealing machine learning models via prediction {APIs}. In *25th USENIX security symposium (USENIX Security 16)*. 601–618.
- [54] Yusuke Uchida, Yuki Nagai, Shigeyuki Sakazawa, and Shin'ichi Satoh. 2017. Embedding watermarks into deep neural networks. In *Proceedings of the 2017 ACM on international conference on multimedia retrieval*. 269–277.
- [55] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).
- [56] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y Zhao. 2019. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)*. IEEE, 707–723.
- [57] Tianhao Wang and Florian Kerschbaum. 2021. Riga: Covert and robust white-box watermarking of deep neural networks. In *Proceedings of the Web Conference 2021*. 993–1004.
- [58] Han Xu, Xiaorui Liu, Yaxin Li, Anil Jain, and Jiliang Tang. 2021. To be robust or to be fair: Towards fairness in adversarial training. In *International conference on machine learning*. PMLR, 11492–11501.
- [59] Yifan Yan, Xudong Pan, Mi Zhang, and Min Yang. 2023. Rethinking White-Box Watermarks on Deep Learning Models under Neural Structural Obfuscation. In *32th USENIX security symposium (USENIX Security 23)*.
- [60] Ziqi Yang, Hung Dang, and Ee-Chien Chang. 2019. Effectiveness of Distillation Attack and Countermeasure on Neural Network Watermarking. *ArXiv abs/1906.06046* (2019).
- [61] Hongxu Yin, Pavlo Molchanov, Jose M Alvarez, Zhizhong Li, Arun Mallya, Derek Hoiem, Niraj K Jha, and Jan Kautz. 2020. Dreaming to distill: Data-free knowledge transfer via deepinversion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8715–8724.
- [62] Youngsik Yoon, Jinhwan Nam, Hyejoeng Yun, Dongwoo Kim, and Jungseul Ok. 2022. Few-Shot Unlearning by Model Inversion. *arXiv preprint arXiv:2205.15567* (2022).
- [63] Jialong Zhang, Zhongshu Gu, Jiyong Jang, Hui Wu, Marc Ph Stoecklin, Heqing Huang, and Ian Molloy. 2018. Protecting intellectual property of deep neural networks with watermarking. In *Proceedings of the 2018 on Asia Conference on Computer and Communications Security*. 159–172.
- [64] Jianpeng Zhang, Yutong Xie, Qi Wu, and Yong Xia. 2019. Medical image classification using synergic deep learning. *Medical image analysis* 54 (2019), 10–19.
- [65] Qi Zhong, Leo Yu Zhang, Shengshan Hu, Longxiang Gao, Jun Zhang, and Yang Xiang. 2022. Attention Distraction: Watermark Removal Through Continual Learning with Selective Forgetting. *2022 IEEE International Conference on Multimedia and Expo (ICME)* (2022), 1–6.

A OMITTED DETAILS OF DEHYDRA

A.1 Classical Natural Image Priors

Specifically, we incorporate two regularization terms $\mathcal{R}_{\ell_2}(\cdot)$ and $\mathcal{R}_{\ell_v}(\cdot)$, as demonstrated in [44], which correspond to ℓ_2 distance and total variation, to steer the optimized samples away from unrealistic ones:

$$\mathcal{R}_{\ell_2}(B) = \sum_{\hat{x}_i \in B} \|\hat{x}_i\|_2^2, \quad (9)$$

$$\mathcal{R}_{\ell_v}(B) = \sum_{\hat{x}_i \in B} \sum_{(j,k)} \sum_{(j',k') \in \delta(j,k)} \|\hat{x}_{i(j,k)} - \hat{x}_{i(j',k')}\|_2^2, \quad (10)$$

where $\delta(j, k)$ indicates a set of pixels adjacent to the pixel at (j, k) in the image \hat{x}_i .

A.2 Splitting Recovered Samples

The detailed pipeline to class-wise split the recovered samples is shown in Alg.1.

Algorithm 1 Split recovered samples.

Input: target model f_w and its l -th layer to analyse, recovered samples $\{B_c\}_*$, saliency ratio β , split ratio γ .

Output: proxy watermark samples $\{B_{c,wmk}\}_*$, proxy normal samples $\{B_{c,nor}\}_*$.

```

1:  $M \leftarrow$  samples' size of each  $B_c$  in  $\{B_c\}_*$ ,
    $D \leftarrow$  neurons' size of the  $l$ -th layer
2:  $\{B_{c,wmk}\}_* \leftarrow \{\}, \{B_{c,nor}\}_* \leftarrow \{\}$ 
3: for batch  $B_c$  in  $\{B_c\}_*$  do
4:   Forward propagate and extract activation of the  $l$ -th layer
      $act = f_{w(l)}(B_c)$ 
     /* Find salient neurons */
5:    $Impt \leftarrow \{\}$ 
6:   for neuron  $j = 0, 1, \dots, D - 1$  do
7:      $Impt[j] \leftarrow \frac{\mu_j}{\sigma_j}$ 
8:   end for
9:    $sorted\_neurons \leftarrow \text{argsort}(Impt)$ 
10:   $salient\_neurons \leftarrow sorted\_neurons[: \beta \times D]$ 
     /* Split recovered batch  $B_c$  */
11:   $Contribution \leftarrow \{\}$ 
12:  for sample  $\hat{x}_i$  in  $B_c$  do
13:     $Contribution[i] \leftarrow \text{sum}(act[i, salient\_neurons])$ 
14:  end for
15:   $sorted\_samples \leftarrow \text{argsort}(Contribution)$ 
16:   $proxy\_wmk \leftarrow sorted\_samples[\gamma \times M :],$ 
      $proxy\_nor \leftarrow sorted\_samples[: \gamma \times M]$ 
17:   $B_{c,wmk} \leftarrow B_c[proxy\_wmk],$ 
      $B_{c,nor} \leftarrow B_c[proxy\_nor]$ 
18:  Append split data  $B_{c,wmk}$  and  $B_{c,nor}$  to  $\{B_{c,wmk}\}_*$  and
      $\{B_{c,nor}\}_*$  respectively
19: end for
20: return  $\{B_{c,wmk}\}_*, \{B_{c,nor}\}_*$ 

```

B OMITTED DETAILS IN EXPERIMENTS

B.1 More Backgrounds on Existing Black-box Model Watermarking Schemes

Our evaluation covers ten existing black-box DNN watermarks, which are categorized as *fixed-class watermarks* (i.e., all the watermark data are paired with the identical target class), and *non-fixed-class watermarks* (i.e., each watermark data is paired with its own target class).

- *Fixed-class watermarks:* (1) **Content**, (2) **Noise**, (3) **Unrelated** [63]. *Content* and *Noise* watermarks draw normal samples from a source class and apply watermark patterns, which are text shapes and random noises respectively, as watermark data. *Unrelated* watermark draws samples from another dataset. All these selected watermark data are annotated with one fixed target class. (4) **Piracy** [35] modifies the values of certain pixels in an image to extreme maximum or minimum values as the trigger, which is then mapped to a fixed class. The position of these pixels and the target class are determined by a signature provided by the user. (5) **EWE** [27] selects samples from a source class, adds square patterns and annotates with a fixed target class. During training, EWE enforces the model to learn entangled features between watermark samples and normal samples of the target class.
- *Non-fixed-class watermarks:* (6) **Adi** [2] leverages a set of abstract out-of-distribution images as watermark data and annotates them randomly. (7) **AFS** [30] generates adversarial examples near the classification boundary of a normally trained model and annotates them with the ground-truth label. The model is then finetuned with these samples added to the training set. (8) **EW** [45] selects normal samples and annotates them with random labels except the ground-truth. During embedding, a normally trained model is applied with an exponential weight operator layer-wise, and then finetuned with watermark samples added to the training set. (9) **Blind** [37] takes ordinary samples and exclusive logos as inputs, uses an encoder to generate watermarks that are almost indistinguishable from the original samples, and uses a discriminator to assist in training. All watermarks are labeled as a fixed class.. (10) **Mark** [20] modifies a portion of the dataset based on the user's signature. The pixel values at specific positions of each image are modified to extreme values according to the signature, and the label of each image is mapped to one of the remaining labels in a specific way based on the signature..

B.2 Implementations of the Target Watermarks

Our implementations mainly referred to some recent open-source watermark framework [33, 42], and strictly followed the specifications and hyperparameters in the original papers. For *Piracy* [35] and *Blind* [37], we set the size for watermark data as 1% of the clean training set, following their papers. For other black-box watermarks, the size of the watermark data is set to 100.

B.3 Performance of Watermarked Models

The performance of the well-trained DNN models protected by ten different black-box watermarks on three benchmark datasets is shown in Table 12.

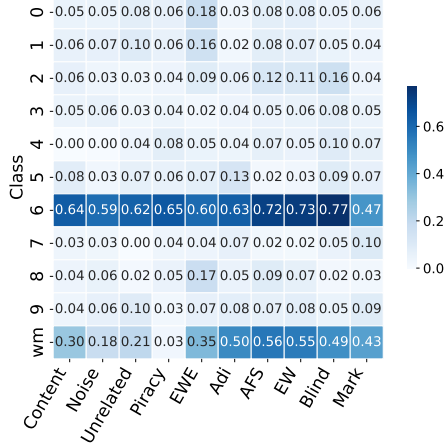


Figure 8: Class-wise intersection ratio between the salient neurons of the real samples and the recovered samples.

Table 12: The performance of the target watermarked models on three datasets. The threshold is estimated on 20 null models as suggested by Lukas et al. [42]

Watermarks	MNIST			CIFAR-10			CIFAR-100		
	Clean / WM	Thre.		Clean / WM	Thre.		Clean / WM	Thre.	
Content	99.08 / 100	0.0083		93.77 / 100	0.1044		76.27 / 100	0.0244	
Noise	99 / 100	0.1057		94.13 / 100	0.4697		76.67 / 100	0.0151	
Unrelated	98.94 / 100	0.088		93.65 / 100	0.0931		75.93 / 100	0.0305	
Piracy	99.22 / 100	0.2807		93.31 / 100	0.1966		73.69 / 100	0.1672	
EWE	98.86 / 100	0.0025		93.97 / 100	0.045		67.73 / 100	0.0059	
Adi	98.82 / 100	0.149		93.91 / 100	0.1425		75.56 / 100	0.0176	
AFS	98.73 / 100	0.1506		92.96 / 100	0.2103		65.76 / 100	0.0813	
EW	98.55 / 100	0.0025		92.89 / 100	0.1335		75.98 / 100	0.0088	
Blind	98.64 / 100	0.1432		91.18 / 100	0.1452		67.22 / 100	0.0266	
Mark	98.57 / 100	0.0025		93.97 / 100	0.1381		65.25 / 100	0.0125	

B.4 Details of Baseline Removal Attacks

Here we provide some implementation details of the baseline removal attacks. For the *Pruning* attack, we set the weight prune ratio to 80%. For the *Fine-pruning* attack, we set the neuron prune ratio to 20%. For the *Finetuning* attack, we start with a learning rate of 0.05 as in [11], followed by a gradually decreasing learning rate schedule. For *Regularization* and *Distraction*, we strictly follow the original papers. For *Laundering*, we follow the main method in [3] and implement it with a certain watermark knowledge. We use the original labels for watermarks using non-OOD data, and least-likely labels for watermarks using OOD data.

B.5 Details of the Compared Recovering Methods

During ablation studies in Section 6.4.2, to better control other variables and perform a fair comparison on the recovering methods, we assume the attacker has got the ground-truth knowledge of the target label settings (i.e., fixed-class with target class 6 for *Content* and *EWE*, and non-fixed for *Adi*) and attempts to perform an adaptive removal attack.

To further ensure the comparison fairness, we implement the improved DEHYDRA with the splitting algorithm omitted. For a

fixed-class watermark, we specifically unlearn the recovered samples corresponding to its target class (or the poisoned samples applied with a trigger towards the target class). For a non-fixed-class watermark, we unlearn the recovered samples from all classes (or the poisoned samples applied with different triggers towards all classes).

C MORE EXPERIMENT RESULTS

C.1 Salient Neurons Analysis of Recovered Samples

In this section, we follow the settings in Section 5.2 and validate the hypothesis of normal data dominance on all the ten black-box watermarks. Concretely, for the ten batches of normal samples X_0, X_1, \dots, X_9 and one batch of watermark samples X_{wm} , we capture the salient neurons of each batch, e.g. $SN(X_i)$, and calculate its intersection ratio w.r.t. those of the recovered samples $SN(B_6)$. Note that for each watermarked model, the sum of intersection ratios of all the classes could be larger than 1, since different classes may share a few common salient neurons.

As is shown in Figure 8, both the fixed-class and non-fixed-class watermarks have witnessed the noisy activations, as the recovered samples would activate salient neurons of samples from classes other than 6. Moreover, $|SN(B_6)| \cap |SN(X_6)|$ is always larger than $|SN(B_6)| \cap |SN(X_{wm})|$, implying that the most consistent activations are correlated with the normal samples of class 6. Note that for non-fixed-class watermarks, the gap still holds but is less evident, possibly due to their more complex input-output pairing relation, which leaves the watermark data more entangled with the normal data.

C.2 Attack Results on MNIST

Note that MNIST is a relatively simple task, and the LeNet-5 [31] structure we select has no batch normalization layers. In this case, we simply set the term $\mathcal{R}_{bn}(B) = 0$ during watermark recovering. Both the basic and improved DEHYDRA still perform well, illustrating the generality of our methods.

We select 1000 samples for the in-distribution setting, and 2000 samples from EMNIST [13] for the transfer setting. The attack results on different data settings are shown in Table 13, Table 14, and Table 15.

C.3 Attack Results on CIFAR-100

For the CIFAR-100 dataset, we select 1000 samples for the in-distribution setting, and 2000 samples from CIFAR-10 for the transfer setting. The attack results on different data settings are shown in Table 16, Table 17, and Table 18.

D PROOF OF WATERMARK SMOOTHNESS DISCREPANCY

D.1 Problem Definition and Theorems

We formulate the *watermark smoothness discrepancy* hypothesis in Theorem D.1 on mixture Gaussian distributions w.r.t. a linear model $f(x) = \text{sign}(\langle w, x \rangle + b)$. These settings are widely used in previous analytical studies for trustworthy machine learning due to the intractable complexity of analyzing DNNs [43, 58].

Table 13: Comparison of removal attacks against black-box watermarks under the in-distribution setting on MNIST. x / y denotes the clean accuracy / rescaled watermark accuracy.

Attacks	Content	Noise	Unrelated	Piracy	EWE	Adi	AFS	EW	Blind	Mark
None	99.1 / 100.0	99.0 / 100.0	98.9 / 100.0	99.2 / 100.0	98.9 / 100.0	98.8 / 100.0	98.7 / 100.0	98.5 / 100.0	98.6 / 100.0	98.6 / 100.0
Fine-pruning	96.6 / 53.6	97.0 / 47.4	97.7 / 100.0	98.4 / 75.0	97.8 / 100.0	96.0 / 52.4	95.6 / 50.0	97.4 / 54.4	97.3 / 49.8	96.0 / 53.4
Finetuning	97.7 / 60.7	97.2 / 45.2	97.4 / 99.5	99.1 / 100.0	97.4 / 87.0	95.3 / 48.3	97.3 / 58.8	95.3 / 54.4	95.9 / 49.2	95.8 / 52.9
Regularization	97.4 / 50.6	97.6 / 45.8	97.2 / 53.9	97.3 / 41.6	97.7 / 51.4	97.4 / 55.3	97.0 / 50.6	95.7 / 54.4	95.9 / 51.0	96.0 / 53.9
Distraction	90.4 / 52.6	87.3 / 44.6	90.8 / 45.7	99.2 / 100.0	95.0 / 49.9	97.9 / 55.9	85.8 / 50.6	85.5 / 55.9	94.4 / 51.6	93.9 / 54.4
Laundering	97.0 / 51.6	97.7 / 46.9	96.0 / 45.2	98.6 / 36.7	96.0 / 53.4	95.4 / 48.3	96.6 / 50.0	95.8 / 55.4	96.9 / 49.8	96.2 / 53.4
DEHYDRA _{Basic}	96.7 / 50.6	96.7 / 44.6	93.8 / 45.2	96.3 / 37.4	97.1 / 49.9	96.7 / 49.5	96.9 / 46.4	95.9 / 53.9	96.1 / 49.2	96.9 / 52.9
DEHYDRA _{Improved}	98.7 / 52.1	99.0 / 45.8	98.7 / 45.2	98.8 / 36.0	98.3 / 51.4	97.2 / 46.5	97.0 / 47.0	96.4 / 53.9	96.5 / 49.8	97.1 / 52.9

Table 14: Comparison of removal attacks against black-box watermarks under the transfer setting on MNIST.

Attacks	Content	Noise	Unrelated	Piracy	EWE	Adi	AFS	EW	Blind	Mark
None	99.1 / 100.0	99.0 / 100.0	98.9 / 100.0	99.2 / 100.0	98.9 / 100.0	98.8 / 100.0	98.7 / 100.0	98.5 / 100.0	98.6 / 100.0	98.6 / 100.0
Regularization	73.1 / 75.8	61.0 / 60.9	65.5 / 99.5	67.8 / 91.0	77.8 / 49.9	72.0 / 58.9	62.6 / 58.2	72.2 / 59.9	80.3 / 53.9	74.7 / 57.4
Distraction	85.9 / 52.1	91.8 / 44.1	96.2 / 81.4	99.2 / 100.0	98.9 / 100.0	95.9 / 61.8	87.5 / 60.0	88.2 / 53.9	98.5 / 90.1	94.8 / 62.4
DEHYDRA _{Basic}	83.2 / 54.1	96.3 / 46.3	95.4 / 51.2	96.9 / 42.3	97.0 / 53.4	94.8 / 61.2	90.9 / 55.3	90.3 / 75.4	91.8 / 64.4	90.8 / 84.0
DEHYDRA _{Improved}	96.8 / 53.1	98.6 / 44.6	97.3 / 45.2	98.1 / 36.0	97.9 / 49.9	95.8 / 60.6	90.6 / 51.1	92.3 / 62.4	91.9 / 62.7	91.1 / 61.9

Table 15: Comparison of removal attacks against black-box watermarks under the data-free setting on MNIST.

Attacks	Content	Noise	Unrelated	Piracy	EWE
None	99.1 / 100.0	99.0 / 100.0	98.9 / 100.0	99.2 / 100.0	98.9 / 100.0
Pruning	98.9 / 100.0	98.3 / 100.0	98.8 / 100.0	99.1 / 100.0	98.5 / 99.5
DEHYDRA (Improved)	93.2 / 49.6	97.8 / 44.1	89.3 / 47.9	97.3 / 36.0	94.5 / 49.9

Concretely, suppose a balanced binary classification dataset $\mathcal{D}_{nor} = (\mathbf{x}, y) \in \mathbb{R}^d \times \{\pm 1\}$, with underlying distributions $\mathcal{N}(\mu_+, \sigma^2 I)$ and $\mathcal{N}(\mu_-, \sigma^2 I)$ for the two classes. The watermark data, $\mathcal{D}_{wm} = \mathcal{N}(\mu_{wm}, \sigma_{wm}^2 I) \times \{+1\}$, have a fixed target class “+1” and are sampled with ratio p compared to the normal dataset during training. Also, we define the input smoothness S_I and parameter smoothness S_p to be the performance of model under input-level and parameter-level Gaussian perturbations respectively. Then we have the following theorem:

THEOREM D.1. Consider the optimal linear model $f^*(x) = \text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*)$ obtained by minimizing the misclassification risk on $\mathcal{D}_{nor} \cup p\mathcal{D}_{wm}$. If the watermark data satisfies (1) $\frac{p}{0.5+p}(\sigma_{wm}^2 - \sigma^2) + \frac{0.5p}{(0.5+p)^2}(\mu_{wm} - \mu_+)^2 > 0$ and (2) $p > 0.5\sqrt{\frac{0.5}{0.5+p} + \frac{p}{0.5+p}\frac{\sigma_{wm}^2}{\sigma^2} + \frac{0.5p}{(0.5+p)^2}\frac{(\mu_{wm} - \mu_+)^2}{\sigma^2}} - 0.5$, then we have $S_{I,+1}(f^*) > S_{I,-1}(f^*)$ and $S_{p,+1}(f^*) > S_{p,-1}(f^*)$.

Remark. The theorem intuitively indicates that, for a fixed-class watermark, the target class will exhibit a higher smoothness (both input-level and parameter-level) compared to other classes if the watermark samples are (1) diverse internally and different from the normal ones in the target class, and (2) sampled much more frequently. During proof, we also find the additive input-level and parameter-level noise would further enlarge the class-wise discrepancies.

D.2 Proof Sketch

We provide the proof sketch for the watermark smoothness discrepancy here. The complete proof is shown in the next section.

PROOF SKETCH. We first consider a simpler mixture Gaussian distribution, where $x \sim \mathcal{N}(1, \sigma_+^2 I)$ for class “+1” with prior probability α and $x \sim \mathcal{N}(-1, \sigma_-^2 I)$ for class “-1” with prior probability $1 - \alpha$. Under this setting, the parameter of the optimal model $f^*(x) = \text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*)$ can be derived by solving the equation where the derivative of the misclassification risk equals zero, i.e. $\frac{dR(f^*)}{df^*} = 0$.

Next, we prove two lemmas on f^* ’s class-wise discrepancies for input smoothness and parameter smoothness, respectively. Specifically, for input smoothness, if $\sigma_+ > \sigma_-$ and $\frac{\alpha\sigma_-}{(1-\alpha)\sigma_+} > 1$, then we have $S_{I,+1}(f^*) > S_{I,-1}(f^*)$. This means that if class “+1” has a larger variance and is sampled more heavily, then f^* should perform better on class “+1” under input-level perturbations, compared with that on class “-1”. A similar conclusion holds for parameter smoothness. During the derivation, we also find that the variance of the additive perturbations is positively correlated with the class-wise discrepancies.

Finally, for a watermarked model trained jointly on $\mathcal{D}_{nor} \cup p\mathcal{D}_{wm}$, we show that the training data belonging to class “+1”, including $\mathcal{N}(\mu_+, \sigma^2 I)$ and $\mathcal{N}(\mu_{wm}, \sigma_{wm}^2 I)$, actually forms a new mixture Gaussian distribution. Therefore, this problem can be easily reduced to the lemmas’ settings above, by translating, rotating and scaling the coordinate system. Substituting the concrete parameters of watermark data into the conditions $\sigma_+ > \sigma_-$ and $\frac{\alpha\sigma_-}{(1-\alpha)\sigma_+} > 1$ yields the final conditions for watermark smoothness discrepancy as in Theorem D.1. \square

D.3 Complete Proof

We first consider a binary classification problem on $\mathcal{D} = (\mathbf{x}, y) \in \mathbb{R}^d \times \{\pm 1\}$ with classes “+1” and “-1”, each following a Gaussian distribution. Without loss of generality, we assume the distribution of \mathcal{D} as:

$$y = \begin{cases} +1, p = \alpha \\ -1, p = 1 - \alpha \end{cases}, \quad \mathbf{x} \sim \begin{cases} \mathcal{N}(1, \sigma_+^2 I) & \text{if } y = +1 \\ \mathcal{N}(-1, \sigma_-^2 I) & \text{if } y = -1 \end{cases} \quad (11)$$

Table 16: Comparison of removal attacks against black-box watermarks under the in-distribution setting on CIFAR-100.

Attacks	Content	Noise	Unrelated	Piracy	EWE	Adi	AFS	EW	Blind	Mark
None	76.3 / 100.0	76.7 / 100.0	75.9 / 100.0	73.7 / 100.0	67.7 / 100.0	75.6 / 100.0	65.8 / 100.0	76.0 / 100.0	67.2 / 100.0	65.2 / 100.0
Fine-pruning	56.3 / 53.9	58.6 / 54.3	58.4 / 48.4	7.6 / 40.0	63.3 / 95.5	57.5 / 56.2	51.2 / 70.6	52.6 / 61.7	10.6 / 49.1	51.3 / 61.5
Finetuning	67.3 / 100.0	68.5 / 100.0	67.9 / 100.0	12.5 / 40.0	60.1 / 99.0	69.5 / 99.5	63.2 / 98.9	68.5 / 96.5	31.2 / 59.4	63.4 / 98.0
Regularization	58.2 / 100.0	59.2 / 100.0	60.4 / 100.0	60.5 / 40.0	25.9 / 49.7	57.1 / 100.0	31.9 / 69.0	59.5 / 96.5	36.1 / 70.7	34.7 / 52.4
Distraction	56.1 / 48.7	50.1 / 49.2	60.9 / 48.4	58.7 / 46.0	64.7 / 100.0	59.8 / 51.1	63.7 / 99.5	56.6 / 50.6	27.0 / 49.1	64.4 / 85.8
Laundering	67.6 / 49.8	68.0 / 49.2	67.8 / 48.4	46.7 / 40.6	47.3 / 49.7	58.9 / 71.0	60.1 / 96.7	66.9 / 93.9	33.3 / 49.1	62.1 / 91.9
DEHYDRABasic	72.6 / 48.7	63.2 / 57.9	63.0 / 55.6	65.8 / 44.8	64.5 / 49.7	62.1 / 57.2	65.2 / 57.5	61.2 / 57.1	54.5 / 58.4	64.0 / 60.5
DEHYDRAImproved	75.1 / 48.7	73.9 / 50.2	73.1 / 53.1	71.0 / 41.8	67.1 / 49.7	64.8 / 60.3	65.4 / 54.8	62.8 / 59.6	65.1 / 55.8	64.5 / 58.5

Table 17: Comparison of removal attacks against black-box watermarks under the transfer setting on CIFAR-100.

Attacks	Content	Noise	Unrelated	Piracy	EWE	Adi	AFS	EW	Blind	Mark
None	76.3 / 100.0	76.7 / 100.0	75.9 / 100.0	73.7 / 100.0	67.7 / 100.0	75.6 / 100.0	65.8 / 100.0	76.0 / 100.0	67.2 / 100.0	65.2 / 100.0
Regularization	59.3 / 100.0	53.9 / 100.0	52.6 / 100.0	48.4 / 100.0	42.3 / 49.7	60.3 / 100.0	52.2 / 90.2	57.1 / 95.5	43.8 / 88.2	49.3 / 52.4
Distraction	76.3 / 90.3	76.4 / 100.0	76.0 / 99.5	72.3 / 52.6	64.9 / 100.0	56.3 / 50.1	63.1 / 97.8	59.4 / 57.6	66.7 / 98.5	63.8 / 77.7
DEHYDRABasic	72.7 / 48.7	56.6 / 51.3	62.3 / 58.2	64.1 / 50.2	56.2 / 59.8	59.1 / 58.3	58.4 / 54.8	59.5 / 59.1	55.7 / 54.8	55.4 / 60.0
DEHYDRAImproved	74.0 / 48.7	75.2 / 49.7	70.0 / 58.2	71.0 / 44.8	62.2 / 59.3	55.1 / 52.7	58.8 / 49.9	61.9 / 58.6	56.1 / 54.3	61.1 / 58.5

Table 18: Comparison of removal attacks against black-box watermarks under the data-free setting on CIFAR-100.

Attacks	Content	Noise	Unrelated	Piracy	EWE
None	76.3 / 100.0	76.7 / 100.0	75.9 / 100.0	73.7 / 100.0	67.7 / 100.0
Pruning	66.4 / 100.0	66.9 / 100.0	58.7 / 100.0	26.7 / 40.0	60.2 / 100.0
DEHYDRA (Improved)	70.0 / 50.8	72.2 / 51.8	61.1 / 58.7	71.1 / 41.8	57.5 / 60.3

where α is the prior probability for class “+1”, $\mathbf{1} = (1, 1, \dots, 1)$ is a d -dim vector and I is a d -dim identity matrix. This simple setting will be adapted to a more complex one in Section D.3.4, with additional watermark data incorporated during training.

We firstly derive the parameter form of the optimal model obtained during training, and analyze the class-wise smoothness discrepancies of the model during testing subsequently.

D.3.1 Optimal Model. Denote $f^*(x) = \text{sign}(\langle \mathbf{w}^*, x \rangle + b^*)$ to be the optimal model. Due to the assumed Gaussian distributions with origin symmetry, we can prove $w_1^* = w_2^* = \dots = w_d^* > 0$ for \mathbf{w}^* by contradiction. For simplicity, we can also define $\eta^* = \frac{b^*}{w_1^*}$, which is the parameter to solve for f^* .

Based on the properties of Gaussian distributions, the empirical risk to minimize for class “+1” is:

$$\begin{aligned}
 R_{+1}(f^*) &= \mathbb{P}\left(\sum_{i=1}^d w_i^* x_i + b^* < 0\right) = \mathbb{P}\left(w_1^* \sum_{i=1}^d x_i + b^* < 0\right) \\
 &= \Phi\left(-\frac{b^* + dw_1^*}{\sqrt{dw_1^* \sigma_+}}\right) = \Phi\left(-\frac{\eta^* + d}{\sqrt{d\sigma_+}}\right),
 \end{aligned} \tag{12}$$

and the empirical risk for class “-1” is:

$$\begin{aligned}
 R_{-1}(f^*) &= \mathbb{P}\left(\sum_{i=1}^d w_i^* x_i + b^* > 0\right) = \mathbb{P}\left(w_1^* \sum_{i=1}^d x_i + b^* > 0\right) \\
 &= 1 - \Phi\left(-\frac{b^* + dw_1^*}{\sqrt{dw_1^* \sigma_-}}\right) = \Phi\left(\frac{\eta^* - d}{\sqrt{d\sigma_-}}\right).
 \end{aligned} \tag{13}$$

Therefore, the total objective function for f^* is:

$$R(f^*) = \alpha \Phi\left(-\frac{\eta^* + d}{\sqrt{d\sigma_+}}\right) + (1 - \alpha) \Phi\left(\frac{\eta^* - d}{\sqrt{d\sigma_-}}\right). \tag{14}$$

We can find the optimal η^* minimizing the above objective function by taking $\frac{dR(f^*)}{d\eta^*} = 0$, which can be expanded as:

$$\begin{aligned}
 &\alpha \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{\eta^* + d}{\sqrt{d\sigma_+}}\right)^2\right) \frac{-1}{\sqrt{d\sigma_+}} + \\
 &(1 - \alpha) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{\eta^* - d}{\sqrt{d\sigma_-}}\right)^2\right) \frac{1}{\sqrt{d\sigma_-}} = 0.
 \end{aligned} \tag{15}$$

The solution to this equation is:

$$\eta^* = d \frac{\sigma_+^2 + \sigma_-^2 - 2\sigma_+ \sigma_- \sqrt{1 - \frac{K}{2d}(\sigma_+^2 - \sigma_-^2)}}{\sigma_+^2 - \sigma_-^2}, \tag{16}$$

where $K = \log(\frac{\alpha\sigma_-}{(1-\alpha)\sigma_+})$. Note that we drop another solution to Equation 15 since its decision boundary does not lie between the two Gaussians. In this way, we derive the parameter η^* of the optimal model f^* for \mathcal{D} .

D.3.2 Input Smoothness Discrepancy. We first study the input smoothness $S_{I,+1}(f^*)$ and $S_{I,-1}(f^*)$ of the two classes, i.e. the performance of model on the class samples under input-level Gaussian perturbations $\mathcal{N}(\mathbf{0}, \sigma_I^2 I)$.

LEMMA D.2. Consider the linear model $f^*(x) = \text{sign}(\langle \mathbf{w}^*, x \rangle + b^*)$ obtained by minimizing Equation 14 on mixture Gaussian distribution \mathcal{D} . If $\sigma_+ > \sigma_-$ and $\frac{\alpha\sigma_-}{(1-\alpha)\sigma_+} > 1$, then we have $S_{I,+1}(f^*) > S_{I,-1}(f^*)$.

PROOF. Denote ϵ_I to be the noise added to each input element. The input smoothness of class “+1” and “-1” can be expressed as follows:

$$S_{I,+1}(f^*) = \mathbb{P}\left(\sum_{i=1}^d w_i^* (x_i + \epsilon_I) + b^* > 0\right) = \Phi\left(\frac{\eta^* + d}{\sqrt{d(\sigma_+^2 + \sigma_I^2)}}\right), \tag{17}$$

$$S_{I,-1}(f^*) = \mathbb{P} \left(\sum_{i=1}^d w_i^*(x_i + \epsilon_I) + b^* < 0 \right) = \Phi \left(\frac{-\eta^* + d}{\sqrt{d(\sigma_-^2 + \sigma_I^2)}} \right). \quad (18)$$

Here we can find that

$$\begin{aligned} \frac{\Phi^{-1}(S_{I,+1}(f^*))}{\sqrt{d}} - \frac{\Phi^{-1}(S_{I,-1}(f^*))}{\sqrt{d}} &= \frac{1 + \frac{\eta^*}{d}}{\sqrt{\sigma_+^2 + \sigma_I^2}} - \frac{1 - \frac{\eta^*}{d}}{\sqrt{\sigma_-^2 + \sigma_I^2}} \\ &= \frac{\sqrt{\sigma_-^2 + \sigma_I^2} \cdot (2\sigma_+^2 - 2\sigma_+\sigma_- - \sqrt{1 - \frac{K}{2d}(\sigma_+^2 - \sigma_-^2)})}{\sqrt{\sigma_+^2 + \sigma_I^2} \cdot \sqrt{\sigma_-^2 + \sigma_I^2} \cdot (\sigma_+^2 - \sigma_-^2)} \\ &\quad - \frac{\sqrt{\sigma_+^2 + \sigma_I^2} \cdot (-2\sigma_-^2 + 2\sigma_+\sigma_- - \sqrt{1 - \frac{K}{2d}(\sigma_+^2 - \sigma_-^2)})}{\sqrt{\sigma_+^2 + \sigma_I^2} \cdot \sqrt{\sigma_-^2 + \sigma_I^2} \cdot (\sigma_+^2 - \sigma_-^2)} \\ &>_{(1)} \frac{\sqrt{\sigma_-^2 + \sigma_I^2} \cdot (2\sigma_+^2 - 2\sigma_+\sigma_-) - \sqrt{\sigma_+^2 + \sigma_I^2} \cdot (-2\sigma_-^2 + 2\sigma_+\sigma_-)}{\sqrt{\sigma_+^2 + \sigma_I^2} \cdot \sqrt{\sigma_-^2 + \sigma_I^2} \cdot (\sigma_+^2 - \sigma_-^2)} \\ &= \frac{2\sigma_+\sigma_- \left(\sqrt{1 + \frac{\sigma_-^2}{\sigma_+^2}} - \sqrt{1 + \frac{\sigma_+^2}{\sigma_-^2}} \right)}{\sqrt{\sigma_+^2 + \sigma_I^2} \cdot \sqrt{\sigma_-^2 + \sigma_I^2} \cdot (\sigma_+ + \sigma_-)} >_{(2)} 0. \end{aligned} \quad (19)$$

Therefore, we can prove $S_{I,+1}(f^*) > S_{I,-1}(f^*)$.

Note that the first inequality $>_{(1)}$ above exploits $\sqrt{1 - \frac{K}{2d}(\sigma_+^2 - \sigma_-^2)} < 1$ since $\sigma_+ > \sigma_-$ and $K = \log(\frac{\alpha\sigma_-}{(1-\alpha)\sigma_+}) > 0$. The second inequality $>_{(2)}$ holds due to the extra noise $\mathcal{N}(0, \sigma_I^2 I)$, **which further enlarges the class-wise discrepancy**. In the extreme case of $\sigma_I = 0$, this inequality turns to an equality, and the smoothness discrepancy degenerates to natural performance discrepancy. \square

D.3.3 Parameter Smoothness Discrepancy. We also study the parameter smoothness $S_{P,+1}(f^*)$ and $S_{P,-1}(f^*)$ here, i.e. the performance of model on the class samples under parameter-level (weights and biases) Gaussian perturbations $\mathcal{N}(0, \sigma_p^2 I)$. Our finding is similar to Lemma D.2.

LEMMA D.3. Consider the linear model $f^*(x) = \text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*)$ obtained by minimizing Equation 14 on mixture Gaussian distribution \mathcal{D} . If $\sigma_+ > \sigma_-$ and $\frac{\alpha\sigma_-}{(1-\alpha)\sigma_+} > 1$, then we have $S_{P,+1}(f^*) > S_{P,-1}(f^*)$.

PROOF. Since the product of two Gaussian random variables is hard to process, we approximately study the parameter smoothness on the centers of $\mathcal{N}(1, \sigma_+^2 I)$ and $\mathcal{N}(-1, \sigma_-^2 I)$, respectively. Denote ϵ_P to be the noise added to each parameter element. The parameter smoothness of class “+1” and “-1” can be expressed as follows:

$$\begin{aligned} S_{P,+1}(f^*) &= \mathbb{P} \left(\sum_{i=1}^d (w_i^* + \epsilon_P) \cdot 1 + (b^* + \epsilon_P) > 0 \right) \\ &= \Phi \left(\frac{b^* + d w_1^*}{\sqrt{d + 1\sigma_P}} \right) = \Phi \left(\frac{w_1^*(d + \eta^*)}{\sqrt{d + 1\sigma_P}} \right), \end{aligned} \quad (20)$$

$$\begin{aligned} S_{P,-1}(f^*) &= \mathbb{P} \left(\sum_{i=1}^d (w_i^* + \epsilon_P) \cdot (-1) + (b^* + \epsilon_P) < 0 \right) \\ &= \Phi \left(\frac{-b^* + d w_1^*}{\sqrt{d + 1\sigma_P}} \right) = \Phi \left(\frac{w_1^*(d - \eta^*)}{\sqrt{d + 1\sigma_P}} \right). \end{aligned} \quad (21)$$

Note that w_1^* is proved to be positive in Section D.3.1. Since $\sigma_+ > \sigma_-$ and $K = \log(\frac{\alpha\sigma_-}{(1-\alpha)\sigma_+}) > 0$, we have $\sqrt{1 - \frac{K}{2d}(\sigma_+^2 - \sigma_-^2)} < 1$ and hence $\eta^* > d \cdot \frac{\sigma_+^2 + \sigma_-^2 - 2\sigma_+\sigma_- - 1}{\sigma_+^2 - \sigma_-^2} = d \cdot \frac{(\sigma_+ - \sigma_-)^2}{\sigma_+^2 - \sigma_-^2} > 0$. Therefore, we can prove $S_{P,+1}(f^*) > S_{P,-1}(f^*)$. \square

D.3.4 Adaptation to a Watermarked Model. Finally, we study the smoothness property of a linear model with fixed-class watermark (target class is “+1”) in a more realistic setting. For a normal dataset $\mathcal{D}_{\text{nor}} = (\mathbf{x}, y) \in \mathbb{R}^d \times \{\pm 1\}$, suppose the distribution is:

$$y = \begin{cases} +1, p = 0.5 \\ -1, p = 0.5, \end{cases} \quad \mathbf{x} \sim \begin{cases} \mathcal{N}(\boldsymbol{\mu}_+, \sigma^2 I) & \text{if } y = +1 \\ \mathcal{N}(\boldsymbol{\mu}_-, \sigma^2 I) & \text{if } y = -1, \end{cases} \quad (22)$$

and the watermark data, $\mathcal{D}_{\text{wm}} = \mathcal{N}(\boldsymbol{\mu}_{\text{wm}}, \sigma_{\text{wm}}^2 I) \times \{+1\}$, are sampled with ratio p compared to the normal dataset during training.

THEOREM D.4. Consider the linear model $f^*(x) = \text{sign}(\langle \mathbf{w}^*, \mathbf{x} \rangle + b^*)$ obtained by minimizing Equation 14 on $\mathcal{D}_{\text{nor}} \cup p\mathcal{D}_{\text{wm}}$. If the watermark data satisfies $\frac{p}{0.5+p}(\sigma_{\text{wm}}^2 - \sigma^2) + \frac{0.5p}{(0.5+p)^2}(\boldsymbol{\mu}_{\text{wm}} - \boldsymbol{\mu}_+)^2 > 0$

and $p > 0.5 \sqrt{\frac{0.5}{0.5+p} + \frac{p}{0.5+p} \frac{\sigma_{\text{wm}}^2}{\sigma^2} + \frac{0.5p}{(0.5+p)^2} \frac{(\boldsymbol{\mu}_{\text{wm}} - \boldsymbol{\mu}_+)^2}{\sigma^2}} - 0.5$, then we have $S_{I,+1}(f^*) > S_{I,-1}(f^*)$ and $S_{P,+1}(f^*) > S_{P,-1}(f^*)$.

PROOF. The training data belonging to class “+1” forms a mixture Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}_{\text{mix}}, \sigma_{\text{mix}}^2 I)$ with

$$\begin{aligned} \boldsymbol{\mu}_{\text{mix}} &= \frac{0.5}{0.5+p} \boldsymbol{\mu}_+ + \frac{p}{0.5+p} \boldsymbol{\mu}_{\text{wm}}, \\ \sigma_{\text{mix}}^2 &= \frac{0.5}{0.5+p} \sigma^2 + \frac{p}{0.5+p} \sigma_{\text{wm}}^2 + \frac{0.5p}{(0.5+p)^2} (\boldsymbol{\mu}_{\text{wm}} - \boldsymbol{\mu}_+)^2. \end{aligned} \quad (23)$$

The model f^* are jointly trained on $\mathcal{N}(\boldsymbol{\mu}_{\text{mix}}, \sigma_{\text{mix}}^2 I) \times \{+1\}$ and $\mathcal{N}(\boldsymbol{\mu}_-, \sigma^2 I) \times \{-1\}$, with prior probability $\alpha = \frac{0.5+p}{1+p}$ for class “+1”.

We can reduce this problem to the lemmas above by translating, rotating the coordinate system and scaling by $\frac{2\sqrt{d}}{\|\boldsymbol{\mu}_{\text{mix}} - \boldsymbol{\mu}_-\|}$, to project the distribution $\mathcal{N}(\boldsymbol{\mu}_{\text{mix}}, \sigma_{\text{mix}}^2 I)$ to $\mathcal{N}(1, \frac{4d}{(\boldsymbol{\mu}_{\text{mix}} - \boldsymbol{\mu}_-)^2} \sigma_{\text{mix}}^2 I)$ and $\mathcal{N}(\boldsymbol{\mu}_-, \sigma^2 I)$ to $\mathcal{N}(-1, \frac{4d}{(\boldsymbol{\mu}_{\text{mix}} - \boldsymbol{\mu}_-)^2} \sigma^2 I)$.

In this way, the condition $\frac{p}{0.5+p}(\sigma_{\text{wm}}^2 - \sigma^2) + \frac{0.5p}{(0.5+p)^2}(\boldsymbol{\mu}_{\text{wm}} - \boldsymbol{\mu}_+)^2 > 0$ yields

$$\sigma_{\text{mix}}^2 > \sigma^2, \quad (24)$$

and $p > 0.5 \sqrt{\frac{0.5}{0.5+p} + \frac{p}{0.5+p} \frac{\sigma_{\text{wm}}^2}{\sigma^2} + \frac{0.5p}{(0.5+p)^2} \frac{(\boldsymbol{\mu}_{\text{wm}} - \boldsymbol{\mu}_+)^2}{\sigma^2}} - 0.5$ yields

$$\frac{0.5+p}{0.5} > \frac{\sigma_{\text{mix}}}{\sigma}. \quad (25)$$

Finally, we obtain $\frac{2\sqrt{d}}{\|\boldsymbol{\mu}_{\text{mix}} - \boldsymbol{\mu}_-\|} \sigma_{\text{mix}} > \frac{2\sqrt{d}}{\|\boldsymbol{\mu}_{\text{mix}} - \boldsymbol{\mu}_-\|} \sigma$ and

$\frac{\alpha \cdot \frac{2\sqrt{d}}{\|\boldsymbol{\mu}_{\text{mix}} - \boldsymbol{\mu}_-\|} \sigma}{(1-\alpha) \cdot \frac{2\sqrt{d}}{\|\boldsymbol{\mu}_{\text{mix}} - \boldsymbol{\mu}_-\|} \sigma_{\text{mix}}} > 1$. Applying Lemma D.2 and D.3, we can conclude that $S_{I,+1}(f^*) > S_{I,-1}(f^*)$ and $S_{P,+1}(f^*) > S_{P,-1}(f^*)$. \square