

QUEEN: Query Unlearning against Model Extraction

Huajie Chen, Tianqing Zhu*, Lefeng Zhang, Bo Liu, Derui Wang, Wanlei Zhou, and Minhui Xue

Abstract—Model extraction attacks currently pose a non-negligible threat to the security and privacy of deep learning models. By querying the model with a small dataset and using the query results as the ground-truth labels, an adversary can steal a piracy model with performance comparable to the original model. Two key issues that cause the threat are, on the one hand, accurate and unlimited queries can be obtained by the adversary; on the other hand, the adversary can aggregate the query results to train the model step by step. The existing defenses usually employ model watermarking or fingerprinting to protect the ownership. However, these methods cannot proactively prevent the violation from happening. To mitigate the threat, we propose QUEEN (QUeRY unLEarNing) that proactively launches counterattacks on potential model extraction attacks from the very beginning. To limit the potential threat, QUEEN has sensitivity measurement and outputs perturbation that prevents the adversary from training a piracy model with high performance. In sensitivity measurement, QUEEN measures the single query sensitivity by its distance from the center of its cluster in the feature space. To reduce the learning accuracy of attacks, for the highly sensitive query batch, QUEEN applies query unlearning, which is implemented by gradient reverse to perturb the softmax output such that the piracy model will generate reverse gradients to worsen its performance unconsciously. Experiments show that QUEEN outperforms the state-of-the-art defenses against various model extraction attacks with a relatively low cost to the model accuracy. The artifact is publicly available at https://anonymous.4open.science/r/queen_implementation-5408/.

Index Terms—Model extraction attacks, disruption-based defenses, sensitivity measurement, AI security.

I. INTRODUCTION

Having achieved revolutionary breakthroughs in various domains, deep neural networks (DNNs) are currently being employed in diverse areas to solve sophisticated real-world problems. The cost of training a high-performance DNN is non-negligible due to the high cost of the large volume of dataset collection, long training time, hardware consumption, etc. For instance, Microsoft has spent 1 billion dollars on OpenAI to develop a large language deep neural network models [1]. Therefore, the deep learning models are considered valuable intellectual properties to be protected by the model owners (referred to as *defenders* in the following). Nowadays, *Model Extraction Attack* (MEA)[2] is considered one of the most critical threats to DNN properties. In MEA, an adversary can establish a piracy model that has the same functionality and comparable performance as that of the protectee model

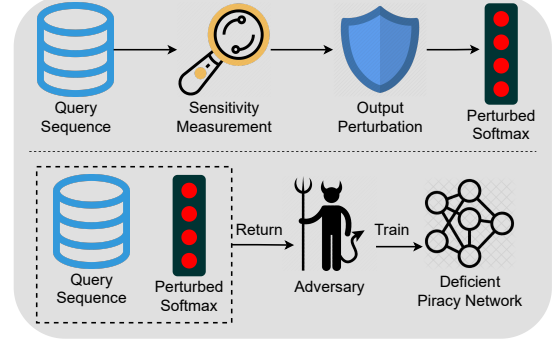


Fig. 1. The overview of QUEEN. For any sequence of queries, QUEEN first measures the sensitivity of each query, and then performs output perturbation accordingly to create perturbed softmax outputs. The piracy network trained with such queries and perturbed outputs will not have comparable performance to the original model.

by sending queries to the original models. This type of attack can lead to serious consequences in terms of copyrights as the adversary can steal from real-world Machine Learning as a Service (MLaaS) models at a very low cost. Meanwhile, MEAs are often used as the initial step for other attacks such as membership inference or model inversion attacks, where the adversary trains shadow models in the manner of MEAs [3], [4]. Hence, we urgently need a series of defenses to counter MEAs, especially in the era of large model.

To mitigate the threat of MEAs, extensive efforts have been made to enhance defense strategies. The existing defense methods can be classified into two classes, passive and active defenses. Passive defenses mainly comprise model watermarking and fingerprinting, where the defender extracts the embedded watermark from the piracy model or uses the fingerprints to get pre-defined output as proof of copyright violation [5], [6], [7], [8], [9]. Despite the effectiveness of the passive defenses, the validation can only be conducted *after the violation*, and the adversary could have already benefited from the piracy model before the crime is discovered. Moreover, if the adversary uses the piracy model only as a proxy to launch further attacks, or the adversary keeps the piracy model for private usage, the defender can no longer protect the protectee model with the passive defenses.

In contrast, active defenses [10], [11], [12] proactively prevent MEAs from happening by jeopardizing the training process with deliberately falsified outputs or lengthening the query time. Active defenses decrease the performance of the piracy models by perturbing the outputs of the protectee model, resulting in a degraded performance of the protectee

*Tianqing Zhu is the corresponding author.

Huajie Chen and Bo Liu are with University of Technology Sydney.

Tianqing Zhu, Lefeng Zhang, and Wanlei Zhou are with City University of Macau.

Derui Wang and Minhui Xue are with CSIRO Data61.

model. Additionally, the majority of active defenses perturb the outputs randomly, which degrades the protectee model's performance significantly. Alternatively, the defender can also extend the query time of the adversary. The current MLaaS provides most services in a black-box setting, where the user can only query the model and receive the query outputs. Hence, we aim to address two questions to counter the MEAs: 1) *Can we detect threats from the queries?* 2) *Can we proactively eliminate the threat before the violation occurs?*

We have explored the possibilities to answer above questions. For the first question, we find that the threat queries show potential trend to cover all classes in the DNN model. In other words, the coverage of those threaten queries is much broader than normal queries series. This enlightens us to estimate the coverage of the query set to detect threats. For the second question, when we identify the threat, it is possible for defender to provide some perturbed query answers that can mislead the shadow models, so it is possible to eliminate the threat before it actually happens.

Due to the above positive rationale, we propose QUEEN (QUERy unEarNing), which proactively launches counterattacks on potential security risks accompanied by sensitivity measurement. As depicted in Figure 1, there are two main components in QUEEN, namely **sensitivity measurement** and **output perturbation**. Given an arbitrary query sequence, the defender estimates the risk of the query by computing the single query sensitivity (SQS) and iteratively updating the cumulative query sensitivity (CQS) class-wise. In output perturbation, if the CQS of the predicted class of the query does not exceed the pre-defined threshold, the defender will perturb the query in the feature space to obfuscate the softmax output. If the CQS exceeds the threshold, the defender launches gradient reverse to make the piracy model generate gradients that worsen the model during training.

Several challenges are identified and resolved during our research: 1) *How to effectively measure the sensitivity of single queries?* For each class in the training dataset, the defender extracts the features of the training samples, and defines the cluster center of each class as the most sensitive point. The distance between the feature of the query and the cluster center is used to estimate the SQS. To estimate the CQS, each query feature is treated as a hypersphere. By computing the ratio of the query features' volume to the volume of the sensitive region, the defender can update the CQS by class query-wise. 2) *How to guarantee that the piracy model produces reverse gradients when trained using the falsified confidence vector?* To launch gradient reverse, the defender must estimate the output of the piracy model given an input. Thus, the defender trains a set of shadow models using subsets of the training dataset to represent the incomplete piracy model. To increase the randomness of the results, a fraction of the shadow models are randomly drawn from the set to produce the average confidence vector for each query. With the estimated piracy confidence vector and that of the protectee model, the defender is able to launch gradient reverse.

Our contributions are summarized as follows.

- We propose a novel counterattack against model extraction named QUEEN that proactively sabotages the MEAs

TABLE I
SUMMARY OF NOTATIONS.

Symbols	Definitions
\mathcal{D}, \mathcal{F}	The problem domain/feature space
D, D'	The original/auxiliary dataset
\mathbf{x}, \mathbf{x}'	The original/auxiliary data sample
$\mathbf{y}, \hat{\mathbf{y}}, \tilde{\mathbf{y}}$	The original/predicted/perturbed confidence vector
y, \hat{y}, \tilde{y}	The original/predicted/perturbed label
f, h	The protectee/piracy model
θ, τ	The parameters of the protectee/piracy model
n, m	The number of samples in D/D'
f^E, f^C	The feature extraction/classification block of f
U, Z	The feature set in \mathcal{F}/\mathbb{R}^2

before the copyright violation happens. QUEEN can be efficiently integrated with generic classification networks without interfering with the training process.

- We have designed novel sensitivity measurement and output perturbation algorithms that can efficiently and precisely prevent MEAs while maintaining the prediction accuracy of the protectee model. The output is selectively perturbed to worsen the piracy model.
- We have conducted extensive experiments to validate the effectiveness of QUEEN, where QUEEN has outperformed the SOTA defenses.

II. PRELIMINARIES AND RELATED WORK

The notations are listed in Table I.

A. Model Extraction Attacks

Given the protectee model $f(\cdot; \theta) : \mathbb{R}^m \mapsto \mathbb{R}^n$ parameterized by θ , the objective of MEAs is to establish a piracy model $h(\cdot; \tau) : \mathbb{R}^m \mapsto \mathbb{R}^n$ parameterized by τ , where h mimics the functionality of f by optimizing

$$\arg \min_{\tau} \mathcal{L}(f(\mathbf{x}; \theta), h(\mathbf{x}; \tau)), \quad (1)$$

where \mathcal{L} and $\mathbf{x} \in \mathbb{R}^m$ respectively denote the loss function and the query sample. The adversary has only black-box access to f , where the adversary has zero knowledge about the architecture, parameters or hyperparameters of f . Normally, the adversary is assumed to have a limited number of unlabelled data samples and computational resources, otherwise the adversary can train the model independently. The dataset owned by the adversary can either be in-distribution or out-of-distribution to the training dataset of f . The auxiliary dataset can be collected online.

Querying the protectee model with natural samples is the most direct attack. However, due to the limited query budget, the adversary tends to select the optimal samples for querying. For instance, semi-supervised learning [13], reinforcement learning [14], and active learning [15] are used for query sampling selection. Additionally, using synthetic data to query the protectee model can also achieve considerable attack accuracy. There are various ways to generate the synthetic samples, such as FGSM [16], C&W attack and feature adversary attack [17]. For example, FGSM uses the gradient to modify the sample

x , such that the modified sample deviates from the predicted class given by h , which is defined as

$$x \leftarrow x - \eta \nabla_x \mathcal{L}(x, h(x; \tau)), \quad (2)$$

where η controls the strength of the modification. Furthermore, even without an auxiliary dataset, the adversary can use generative models to launch data-free MEAs [18], [19], but this requires an enormous query budget.

B. Defenses against Model Extraction Attacks

Generally, the current defenses against MEAs can be divided into passive and active defenses.

Passive defenses mainly include model watermarking and fingerprinting, where the defender can determine whether a suspect model is a piracy model or not by validation. Model watermarking [5], [6], [20], [21], [22] aims to embed watermarks into the protectee model during training, where a set of trigger samples can make the protectee model generate the pre-defined outputs. Model fingerprinting [8], [23], [7], in contrast, does not interfere with the training process, and it creates adversarial samples as fingerprints to make the piracy model and the protectee model share the same outputs when given the fingerprints. In summary, passive defenses allow the defender to claim the copyright of the piracy model by ownership verification. However, this does not prevent the violation from the root.

Active defenses proactively prevent MEAs by perturbing the outputs of the protectee model or intentionally increasing the query time when the query number ascends. Orekondy et al. [24] perturb the output by maximizing the angle deviation between the original and the perturbed gradients. Lee et al. [11] attach an additional layer to the end of the protectee model to produce random softmax output while maintaining the argmax of the output unchanged. Juuti et al. [25] use the Shapiro-Wilk statistics test to distinguish the benign queries from the malicious queries, where the benign samples are believed to fit normal distributions and the malicious do not. Kariyappa et al. [12] train a out-of-distribution (OOD) sample detector to differentiate OOD query samples from normal query samples. The query results of OOD samples are perturbed by a misinformation function. Kariyappa et al. [26] train an ensemble of models to detect OOD query samples and perturb the corresponding outputs. Zhang et al. [27] perturb the query sample to the edge of the decision boundary while keeping the argmax of the output unchanged to produce obfuscated outputs. Dziedzic et al. [10] utilize private aggregation of teacher ensembles (PATE) to measure the privacy risk of each query, and increase the query time when the query number grows by forcing the adversary to solve the proof-of-work (PoW) puzzles.

III. PROBLEM DEFINITION

First, we formally define the protectee model and the piracy model. A problem domain is denoted by $\mathcal{D} \subseteq \mathbb{R}^M$, where each element $\mathbf{x} \in \mathcal{D}$ is labeled by one of N classes. We use $\mathbf{y} \in \mathbb{R}^N$ to denote the one-hot encoded label vector. A deep learning model is a function $f(\cdot; \theta) : \mathbb{R}^M \mapsto \mathbb{R}^N$, parameterized by θ .

A protectee model is a deep learning model $f(\cdot; \theta) : \mathbb{R}^M \mapsto \mathbb{R}^N$ with θ as its parameters trained by a defender who owns a dataset $D \subseteq \mathcal{D}$. D comprises $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$. The defender trains the protectee model f by taking gradient steps to optimize θ on

$$\nabla_{\theta} \frac{1}{n} \sum_{i=1}^n \mathcal{L}\left(f_{\text{sm}}(f(\mathbf{x}_i; \theta)), \mathbf{y}_i\right), \quad (3)$$

where $\mathcal{L}(\cdot, \cdot)$ denotes the loss function such as the cross entropy (CE) loss or Kullback-Leibler divergence (KLDiv) loss; $f_{\text{sm}}(\cdot)$ denotes the softmax function.

A piracy model is another deep learning model $h(\cdot; \tau) : \mathbb{R}^M \mapsto \mathbb{R}^N$ with τ as its parameters established by the adversary by launching model extraction attack on the protectee model f . h shares the same or similar functionalities with f . The adversary collects an auxiliary dataset $D' \subseteq \mathcal{D} \setminus D$. This means that $D' = \{\mathbf{x}'_1, \dots, \mathbf{x}'_m\}$ comes from the same problem domain \mathcal{D} but does not overlap with D . Additionally, $\mathbf{x}' \in D'$ is not labeled by the adversary. The adversary feeds \mathbf{x}' into the protectee model f so as to get the softmax output, i.e., the confidence vector $f(\mathbf{x}')$. Similarly, the adversary trains the piracy model h_{τ} taking gradient steps to optimize τ on

$$\nabla_{\tau} \frac{1}{m} \sum_i^m \mathcal{L}\left(f_{\text{sm}}(h(\mathbf{x}'_i; \tau)), f_{\text{sm}}(f(\mathbf{x}'_i; \theta))\right), \forall \mathbf{x}'_i \in D'. \quad (4)$$

A. Threat Model

We consider the scenarios where not only the label of the prediction, but also the entire confidence vector is of the users' concern. For example, the confidence vector can be further utilized for downstream tasks such as OOD detection [28] or membership inference attack [29]. We define the threat model in terms of the capabilities and limitations of the two parties, namely the adversary and the defender.

The Adversary. We assume that the adversary is capable of:

- *Collecting public data.* The adversary can access public datasets to construct an auxiliary dataset for model extraction. The auxiliary dataset can be either in-distribution or out-of-distribution to the training dataset. For instance, given that the protectee model is trained on MNIST dataset, FEMNIST dataset is in-distribution, and CIFAR-10 dataset is out-of-distribution.
- *Accessing the protectee model.* The adversary can query the protectee model with any input with black-box access. It means that the adversary is limited to sending queries to the protectee model and receiving the predictions from the protectee model only. The predictions are presented in the form of complete confidence vectors.
- *Customizing the training algorithm.* The adversary is able to freely select the network architecture, optimization algorithms, loss functions, etc. However, we assume that the adversary is going to select either the CE or KLDiv loss function, because these two are the most commonly used loss functions in terms of training classifiers.
- *Knowing the specific defense method.* The adversary knows which exact defense method is used, and therefore,

the adversary can launch the most threatening attack on the protectee model.

We assume that the adversary is limited to:

- *Knowledge of the protectee model.* The adversary does not know the model architecture, parameters, or hyperparameters of the protectee model.
- *Access to the training data.* The adversary cannot access the original training samples possessed by the defender.

The Defender. In contrast, the defender is capable of:

- *Completely accessing the protectee model.* The defender has full knowledge about the protectee model including its model architecture, parameters and hyperparameters.
- *Completely accessing the original dataset.* The defender can access and modify any samples in the original dataset.
- *Utilizing the queries.* For any query sent by the user, the defender is allowed to perform any operation on the query, including storing and analyzing the query.

Meanwhile, the defender is limited to:

- *Knowledge of the piracy model.* The defender does not know the model architecture, parameters, or hyperparameters of the piracy model.
- *Knowledge of the specific attack.* The attack method employed by the attacker remains unknown to the defender.

The goal of the adversary is to launch model extraction attack on the protectee model so as to obtain a piracy model that performs comparably to the protectee model. Reversely, the defender aims to differentiate the adversaries from the benign users, and proactively jeopardize any possible threat.

The Adaptive Adversary. We further consider that the adaptive adversary knows about the defense mechanism of QUEEN. Given a perturbed confidence vector $\tilde{\mathbf{y}}$, The adaptive adversary aims to recover the clean prediction confidence vector $\hat{\mathbf{y}} = R(\tilde{\mathbf{y}})$ using a recovery function $R(\cdot)$. Thus, adaptive attacks such as D-DAE [2] and pBayes [30] that establishes $R(\cdot)$ pose the greatest threats to QUEEN. We thus test the performance of QUEEN against these adaptive attacks in the experiment.

B. Concepts of Defense Mechanism

Central Data vs. Peripheral Data. Before designing the method, we try to figure out how the classification works in f . Let f^E and f^C respectively denote the feature extraction block and the classification block in f . Then we have

$$f(\mathbf{x}) = f^C(f^E(\mathbf{x})), \quad (5)$$

where f^E maps \mathbf{x} into a feature space $\mathcal{F} \subseteq \mathbb{R}^O$, and f^C maps $f^E(\mathbf{x})$ into \mathbb{R}^N so as to derive the confidence vector. Given the training dataset D MNIST [31], we use f^E to extract the feature $f^E(\mathbf{x}_i)$, $\forall \mathbf{x}_i \in D$. The features are then projected into a 2D space via t-SNE as depicted in Figure 8 in the Supplemental Materials.

Based on our observation, $f^E(\mathbf{x}_i)$ s are separated into clearly distinguishable clusters by their label y_i , which are the black dots. We then process the test dataset with the same procedure as above, where we use different colors to represent

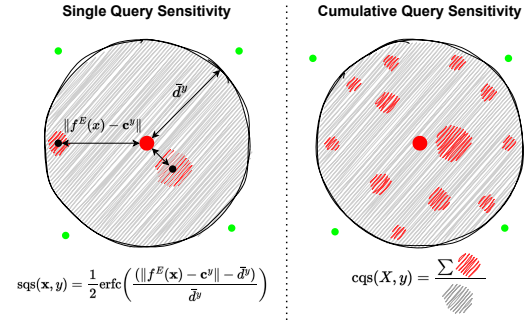


Fig. 2. Sensitivity Measurement. The CQS is the sum of the red/queried area over the gray/sensitive area.

the classes. Two interesting phenomena are observed: 1) the projection of the test data highly overlaps with that of the training data; 2) the misclassified test data points fall in places that are distant from the center of the clusters. Thus, we assume that the center of each cluster in \mathcal{F} represents the most sensitive region in the cluster. Thus, the queries that hit the center of the cluster are the most representative.

To justify our above assumption, we thus conduct a pre-experiment to check whether the central data is more representative than the peripheral data. Here, central data refers to the data whose feature is close to its cluster center, whereas peripheral data refers to the data whose feature is distant from its cluster center. In order to rank the data by their distances to their cluster center, we firstly need to define the cluster center.

Definition 1–Cluster Center: Given a feature cluster $\{f^E(\mathbf{x}_1), \dots, f^E(\mathbf{x}_n)\}$ sharing the same label y , the cluster center is defined as

$$\mathbf{c}^y = \frac{1}{n} \sum_i^n f^E(\mathbf{x}_i), \quad (6)$$

where \mathbf{c}^y denotes the center of the feature cluster labeled by $y = \arg \max(\mathbf{y})$.

The experimental results show that the central data leads to higher performance of the trained classifier network compared to the peripheral data. This suggests that the query whose feature is closer to the cluster center is more sensitive. Section B in the Supplemental Materials shows the details.

Single Query Sensitivity. Given an arbitrary query, the defender aims to know how sensitive it is when the logit produced by the protectee network is honestly returned. Based on the pre-experiment, we know that the probability of a query being more representative than another gets higher, when its feature is closer to the cluster center than the other. Therefore, we define representative probability (RP) to measure the sensitivity of the query.

Definition 2–Single Query Sensitivity: Given an arbitrary query \mathbf{x} labeled by $y = \arg \max \mathbf{y}$, the single query sensitivity $sqs(\mathbf{x}, y)$ is measured by the representative probability $P(\mathbf{x}, y)$, which is computed by

$$\begin{aligned} sqs(\mathbf{x}, y) &= P(\mathbf{x}, y) \\ &= \frac{1}{2} \operatorname{erfc} \left(\frac{(\|f^E(\mathbf{x}) - \mathbf{c}^y\| - \bar{d}^y)}{\bar{d}^y} \right), \end{aligned} \quad (7)$$

where erfc is the complementary error function defined as

$$\text{erfc}(z) = \frac{2}{\sqrt{\pi}} \int_z^\infty e^{-t^2} dt; \quad (8)$$

\bar{d}^y denotes the average distance between \mathbf{c}^y and the features of the training data labeled by y .

As depicted in the left of Figure 2, the black periphery of the circle is the sensitive region determined by \bar{d}^y , and the red dot indicates the cluster center \mathbf{c}^y ; the black dots represent the queries within the sensitive region, whereas the green dots are those out of the region. Intuitively, if the query's feature $f^E(\mathbf{x})$ dynamically gets close to \mathbf{c}^y , $P(\mathbf{x}, y)$ will then increase and be close to 1. Inversely, if $f^E(\mathbf{x})$ gets distant from \mathbf{c}^y , $P(\mathbf{x}, y)$ decreases and eventually end up with 0. Further, the SQS is going to affect the CQS computation by adjusting the radius of the queried space for each query.

Cumulative Query Sensitivity. Given an arbitrary query sequence sharing the same predicted label, the defender aims to quantify the impact caused by honestly returning the confidence vector. As depicted in the right of Figure 2, for all data samples in the training dataset labeled by y , the defender uses f^E to extract the training features. The training features together form a hypersphere dyed gray in the feature space \mathcal{F} , which is called the *sensitive space*. The defender then extracts the query features from the queries and treats each query feature as a small hypersphere dyed red that represents itself and the queries whose features fall around the center of the hypersphere. Together, the query features forms the *queried space*.

By summing up the volume of the query features and dividing it by the volume of the hypersphere formed by the training data, the defender derives the ratio of the queried space to the sensitivity space. We thus define this ratio as the CQS.

Definition 3—Cumulative Query Sensitivity: For a query sequence $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ labeled by $y = \arg \max \mathbf{y}$, the cumulative sensitivity $cqs(X, y)$ is computed by

$$cqs(X, y) = \frac{\sum_{\mathbf{x} \in X} v(sq s(\mathbf{x}, y), r)}{v(sq s(\mathbf{c}^y, y), \bar{d}^y)}, \quad (9)$$

where $v(\cdot, \cdot)$ is the function that computes the volume of a hypersphere; r is a hyperparameter that defines the radius of the hyperspheres of the query features.

If $cqs(\mathbf{x}, y)$ exceeds a pre-defined threshold t , this training class is considered to be threatened by the query sequence. The counterattack is launched to protect the protectee model.

Output Perturbation. If a query sequence categorized into one class is determined to be threatening that class, the defender stops offering the true confidence vector. Instead, the defender proactively jeopardizes the potential possibilities where the queries can support the training of a piracy model. The defender starts to offer falsified confidence vectors to the adversary. Theoretically, the falsified confidence vector is designed to make the piracy model to generate reverse gradient such that the piracy model gets worsened when the falsified confidence vector is used in the training process. We define this objective as follows.

Definition 4—Gradient Reverse: Given an arbitrary query \mathbf{x} , the defender aims to create a falsified confidence vector $\tilde{\mathbf{y}}$ such that

$$\begin{aligned} \nabla_{\tau} \mathcal{L} \left(f_{\text{sm}}(h(\mathbf{x}; \tau)), \tilde{\mathbf{y}} \right) = \\ -\nabla_{\tau} \mathcal{L} \left(f_{\text{sm}}(h(\mathbf{x}; \tau)), f_{\text{sm}}(f(\mathbf{x}; \theta)) \right), \end{aligned} \quad (10)$$

making the gradient completely reverse to the correct direction.

Thereby, the piracy model trained using $(\mathbf{x}, \tilde{\mathbf{y}})$ is crippled, and therefore its performance will be far from comparable to the protectee model.

IV. QUEEN METHOD

A. Overview of QUEEN

There are two main components in our method, namely **sensitivity measurement** and **output perturbation**. As depicted in Figure 3, given the query sequence sent by the user, the defender extracts the query feature with the protectee network. The query feature is then further mapped to a lower-dimensional space for sensitivity measurement. Based on the measurement, different strategies are employed to perturb the softmax output or to return the normal softmax output. The components are further specified as follows.

In sensitivity measurement, a query sensitivity estimation system that estimates the single query sensitivity (SQS) for each query is established. The cumulative query sensitivity (CQS) can thus be derived from the SQS of each query in a given query sequence.

In gradient reverse, given a user and a query sequence, the defender compute the CQS of the query sequence. Then, based on the CQS, the defender determines whether the query sequence is threatening the training dataset class-wise. If there is any potential that the queries of a class could support the training of a piracy model, the defender launches a counterattack by sending falsified confidence vectors. The falsified confidence vectors are designed to worsen the piracy model if they are used in the training process.

B. Sensitivity Analysis

Before making the pre-trained protectee model publicly available, the defender first performs a sensitivity analysis on every class of the training data, so that the defender is able to measure the CQS given any query sequence from the user. The process of sensitivity analysis is demonstrated in Alg. 1, where the objective is to obtain a trained mapping network g , the set of cluster centers \mathbf{c} and the set of average distance \bar{d} .

The defender owns a protectee model f that is trained on the training dataset D . The defender uses the feature extraction layers f^E of f to extract the training features of the training data, which is denoted as $U = \{((u_1, \mathbf{y}_1), \dots, ((u_n, \mathbf{y}_n))\} \subseteq \mathcal{F}$, where $u_i = f^E(\mathbf{x}_i)$. However, it is not practical to compute the volume of the hypersphere in \mathcal{F} , because volume computation in high-dimensional space is complex. Additionally, the defender does not want to modify the parameters θ in f to realize dimension reduction. Therefore, the defender needs a

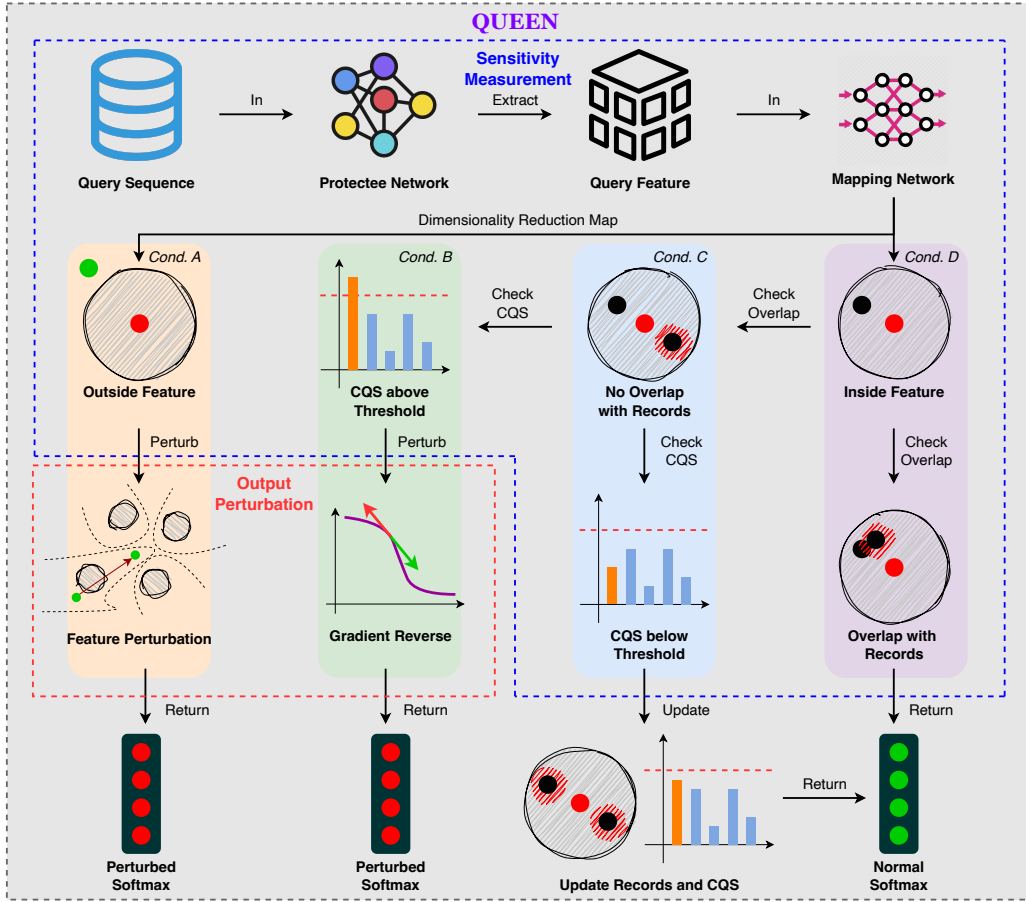


Fig. 3. The workflow of QUEEN. After sensitivity measurement, there are four conditions where the output is either honestly returned or perturbed with gradient reverse or feature perturbation.

Algorithm 1: Sensitivity Analysis

Data: f : the pre-trained protectee model; g : the mapping network; D : the training dataset;
Result: g : the trained mapping network; c : the set of cluster centers; \bar{d} : The set of average distances;
 Use f to extract the training features from D ;
 Train g to map the training features to \mathbb{R}^2 with the supervised contrastive loss till converged;
for each class in D do
 Compute the cluster center;
 Compute the average distance between the cluster center and the 2D features in this class;
 Store the cluster center and the average distance in c and \bar{d} ;
end
return g, c, \bar{d}

fixed mapping network g that maps the features from \mathcal{F} to a low-dimensional space.

In our setting, we define g as a mapping network that maps the features to a two-dimensional (2D) space, denoted as $g : \mathcal{F} \mapsto \mathbb{R}^2$. g can have a very simple architecture, such as several fully-connected layers. Moreover, g is required to make the mapped features preserve the property as that of

the features in \mathcal{F} . That is to say, central features should be more representative than the peripheral features. To achieve this goal, we train g using the supervised contrastive loss[32] that takes the following form.

$$L^{\text{sup}} = \sum_{i \in I} -\log \left\{ \frac{1}{|O(i)|} \sum_{o \in O(i)} \frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_o / \gamma)}{\sum_{a \in A(i)} \exp(\mathbf{z}_i \cdot \mathbf{z}_a / \gamma)} \right\} \quad (11)$$

Here, within a multi-view batch U' randomly sampled from U , $i \in I \equiv \{1, \dots, N\}$ is the index of samples in the batch; \mathbf{z} denotes the mapped 2D feature, namely $\mathbf{z} = g(f^E(\mathbf{x}))$; The operator \cdot denotes the inner product; $A(i) \equiv I \setminus \{i\}$ denotes the set of indices without i ; $O(i) \equiv \{o \in A(i) : \mathbf{y}_o = \mathbf{y}_i\}$ is the set of indices of all other features sharing the same label as that of \mathbf{z}_i ; the operator $|\cdot|$ returns the number of elements in the set; $\gamma \in \mathbb{R}^+$ is a scalar temperature parameter. In short, by minimizing L^{sup} , the defender makes g maps the features in the same class closely, whereas those of different classes are mapped distantly. Thus, the relative position of feature $u_i \in U$ is kept when it is mapped to \mathbf{z}_i . Another feasible solution is to use the hierarchical contrastive loss proposed in [33].

With a trained mapping network g , the defender is able to map entire training feature set F to a 2D feature set denoted by $Z = \{(\mathbf{z}_1, \mathbf{y}_1), \dots, (\mathbf{z}_n, \mathbf{y}_n)\}$. Let $Z^y \subseteq Z$ denote the 2D training feature set of class $y = \arg \max(\mathbf{y})$. The defender

Algorithm 2: Sensitivity Measurement

Data: f : the pre-trained protectee model; D' : the query sequence;

Result: s : the set of CQS of each class;

Use f to get predicted labels for queries in D' ;

Split D' by the predicted labels;

for each split query sequence **do**

 Determine the sensitive region by \bar{d}^y ;

for each query **do**

 Get the 2D feature of the query;

if the feature is within the sensitive region **then**

if no overlapping records **then**

 Compute its SQS;

 Record this feature;

end

end

end

 Get the queried area by summing up the area of each sensitive feature circle;

 Get the CQS of this class by dividing the queried area by the area of the sensitive region;

 Store the CQS in s ;

end

return s

now computes the cluster center c^y in the 2D space for each Z^y as defined in Eq. 6 by

$$c^y = \frac{1}{|Z^y|} \sum_{\mathbf{z}_i \in Z^y} \mathbf{z}_i, \quad (12)$$

Meanwhile, the defender also computes the average distance \bar{d}^y between c^y and $\mathbf{z}_i \in Z^y$ by

$$\bar{d}^y = \frac{1}{|Z^y|} \sum_{\mathbf{z}_i \in Z^y} \|\mathbf{z}_i - c^y\|. \quad (13)$$

C. Sensitivity Measurement

With the mapping network g , the set of cluster center $c = \{\dots, c^y, \dots\}$ and the set of average distances $\bar{d} = \{\dots, \bar{d}^y, \dots\}$, the defender can measure the sensitivity of any given query sequence. As illustrated in Alg. 2, given the auxiliary dataset $D' = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$, the defender gets the set of CQS $s = \{\dots, s^y, \dots\}$, where s^y denotes the CQS of class y .

The defender first gets the predicted label $\hat{y}_i = \arg \max f_{\text{sm}}(f(\mathbf{x}'_i)), \forall \mathbf{x}'_i \in D'$. Next, the defender splits D' by the predicted labels such that $D'^{\hat{y}}$ denotes the subset of D' labeled by \hat{y} . The sensitivity measurement is then conducted class-wise. For each class, the defender determines the sensitive region by setting its radius to $\bar{d}^{\hat{y}}$.

To compute the SQS, the defender gets the 2D feature $\mathbf{z}_i = g(f^E(\mathbf{x}'_i)), \forall \mathbf{x}'_i \in D'^{\hat{y}}$. The defender then needs to determine whether \mathbf{z}_i satisfies the following two conditions:

1. \mathbf{z}_i is within the sensitive region, i.e., $\|\mathbf{z}_i - c^{\hat{y}}\| < \bar{d}^{\hat{y}}$;
2. There is no previously recorded feature overlapping with \mathbf{z}_i , i.e., $\nexists \mathbf{z}_j \in Z^{\hat{y}}, \text{ s.t. } \|\mathbf{z}_i - \mathbf{z}_j\| < r, i \neq j$. Here, r is a hyperparameter denoting the radius of each query circle

representing the features of the query and those nearby it. If the conditions are satisfied, the defender proceeds to compute the SQS of the query by

$$\text{sqs}(\mathbf{z}_i, \hat{y}) = \frac{1}{2} \text{erfc}\left(\frac{\|\mathbf{z}_i - c^{\hat{y}}\| - \bar{d}^{\hat{y}}}{\bar{d}^{\hat{y}}}\right). \quad (14)$$

Eventually, the CQS of the 2D features satisfying the conditions Z with the predicted label \hat{y} is computed by

$$\text{cqs}(Z, \hat{y}) = \frac{r^2}{(\bar{d}^{\hat{y}})^2} \sum_{\mathbf{z}_i \in Z^{\hat{y}}} (\text{sqs}(\mathbf{z}_i, \hat{y}))^2. \quad (15)$$

Intuitively, $\text{cqs}(Z, \hat{y})$ indicates the ratio of the queried area to the sensitive area in class \hat{y} . If $\text{cqs}(Z, \hat{y})$ exceeds a pre-defined threshold, the adversary is determined to be threatening class \hat{y} . Thus, the defender ought to launch a counterattack to prevent further losses.

Eventually, there are four conditions after the sensitivity measurement as depicted in Figure 3:

- *Cond. A: the 2D feature is not within the sensitive region.*
- *Cond. B: the 2D feature is within the sensitive region; the 2D feature overlaps with the records; the CQS exceeds the threshold.*
- *Cond. C: the 2D feature is within the sensitive region; the 2D feature does not overlap with the records; the CQS is below the threshold.*
- *Cond. D: the 2D feature is within the sensitive region; the 2D feature overlaps with the records.*

In *Cond. A/B*, the defender launches feature perturbation/gradient reverse to perturb the softmax outputs, whereas the normal softmax outputs are returned in *Cond. C/D*.

We claim that the data points whose features are not in the sensitive region can also be used in training the piracy model. Our CQS measurement does not include those data points. But those data points are also perturbed by feature perturbation to reduce their contribution to the piracy model training.

D. Output Perturbation

As demonstrated in Alg. 3, the defender proactively defends any potential attacks by conditionally perturbing the confidence vector generated by the protectee model. Given a pre-defined threshold t , if $\text{cqs}(Z, \hat{y}) > t$, the defender starts the gradient reverse counterattack by sending a falsified confidence vector to the adversary to achieve the objective described in Eq. 17. To achieve this objective, the defender needs two main components in gradient reverse, namely *piracy model simulation* and *gradient reverse*.

For the query whose feature is not within the sensitive region, the defender performs *feature perturbation* by moving the feature towards the most distant cluster center. The perturbation stops before the predicted label of the perturbed feature changes, resulting in a confidence vector that worsens the piracy model trains with it.

Piracy Model Simulation. Due to the fact that the defender has zero knowledge about the piracy model, the defender cannot launch gradient reverse because the output of the piracy model is not accessible. To address this problem, the defender

Algorithm 3: Output Perturbation

Data: f : the protectee model; D : the training dataset;
 D' : the query sequence; t : the threshold;
Result: Y : The confidence vectors;
 /* Piracy Model Simulation */
 Randomly split D into subsets evenly;
 Train a set of shadow models of different architectures,
 and each shadow model is only trained on one subset
 without intersection;
 /* Logits Falsification */
for $\mathbf{x} \in D'$ **do**
 | Extract the 2D feature of \mathbf{x} using f and g ;
 | Get the softmax $\hat{\mathbf{y}} \leftarrow f_{\text{sm}}(f(\mathbf{x}))$;
 | Get the label $\hat{y} \leftarrow \arg \max \hat{\mathbf{y}}$;
 | Measure the sensitivity of \mathbf{x} ;
 | **if** *sensitive* **then**
 | | Perform gradient reverse or feature perturbation
 | | based on the conditions to get the falsified
 | | softmax output $\tilde{\mathbf{y}}$;
 | | Store $\tilde{\mathbf{y}}$ in Y ;
 | **else**
 | | Store $\hat{\mathbf{y}}$ in Y ;
 | | Update feature records and CQS accordingly;
 | **end**
end
return Y

needs piracy model simulation, i.e., to train a set of shadow models locally to simulate the behavior of the piracy model. The weights of different models typically converge to the initial stable points within the same optimization problem [34], indicating that the gradient descent directions of the shadow and piracy models are similar. Similarly, the use of shadow models to mimic the piracy model is also illustrated in [24].

The training dataset D is randomly divided into a number of subsets that are used to train shadow models. The shadow models are of various architectures, and each shadow model is paired with one subset without overlapping. In order to simulate the piracy model with randomness, a fixed number of the shadow models are drawn and the average of their logits is used to estimate the piracy model's logit given this input.

Gradient Reverse. Given an arbitrary query $\mathbf{x} \in D'$, the defender first extracts the 2D feature $\mathbf{z} = g(f^E(\mathbf{x}))$, and gets the softmax $\hat{\mathbf{y}} = f_{\text{sm}}(f(\mathbf{x}))$ and the predicted label $\hat{y} = \arg \max \hat{\mathbf{y}}$. If \mathbf{x} is within the sensitive region, and there is no previous record overlapping with it, the query is sensitive. Next, the defender checks whether the CQS of class \hat{y} exceeds the pre-defined threshold t . If yes, the defender starts logits falsification, otherwise it will be treated as a normal query, and $\hat{\mathbf{y}}$ will be honestly returned. The feature of the sensitive query is recorded for overlapping checking.

To perturb the softmax for the sensitive query, the defender randomly draws a subset H' of shadow models from the trained shadow model set H . The falsified logits $\tilde{\mathbf{y}}'$ is derived

Algorithm 4: Feature Perturbation

Data: \mathbf{x} : the query; \mathbf{c} : the cluster centers; f : the protectee model;
Result: $\tilde{\mathbf{y}}$: The perturbed output;
 $u \leftarrow f^E(\mathbf{x})$;
 $\hat{y} \leftarrow \arg \max f_{\text{sm}}(f^C(u))$;
 $y \leftarrow \arg \max \|\mathbf{z} - \mathbf{c}^y\|$;
 $v \leftarrow \frac{\mathbf{z} - \mathbf{c}^y}{\|\mathbf{z} - \mathbf{c}^y\|}$;
while *True* **do**
 | $u' \leftarrow u + \epsilon v$;
 | **if** $\arg \max f_{\text{sm}}(f^C(u')) == \hat{y}$ **then**
 | | $u \leftarrow u'$;
 | **else**
 | | Break;
 | **end**
end
 $\tilde{\mathbf{y}} \leftarrow f_{\text{sm}}(f^C(u))$;
return $\tilde{\mathbf{y}}$

by

$$\hat{\mathbf{y}}' = \frac{1}{|H'|} \sum_{h \in H'} f_{\text{sm}}(h(\mathbf{x})), \quad (16)$$

which is the mean of the logits from the shadow models. With $\hat{\mathbf{y}}$ and $\hat{\mathbf{y}}'$, the defender can now create a falsified logit $\tilde{\mathbf{y}}$ that fulfills the objective of gradient reverse by

$$\tilde{\mathbf{y}} = 2\hat{\mathbf{y}}' - \hat{\mathbf{y}}. \quad (17)$$

In practice, we need to solve an optimization problem so as to get a valid $\tilde{\mathbf{y}}$, because $\tilde{\mathbf{y}}$ derived through Eq. 17 is often not a valid softmax output that can be easily detected by the adversary. The optimization problem is formulated as follows:

$$\begin{aligned} & \max_{\tilde{\mathbf{y}}} \text{cossim}(\tilde{\mathbf{y}}, 2\hat{\mathbf{y}}' - \hat{\mathbf{y}}), \\ & \text{subject to } \sum_{\tilde{y} \in \tilde{\mathbf{y}}} \tilde{y} = 1, \\ & \tilde{y} \geq 0, \forall \tilde{y} \in \tilde{\mathbf{y}}. \end{aligned} \quad (18)$$

After the optimization, $\tilde{\mathbf{y}}$ becomes a valid softmax output where each element in it is greater than 0, and the sum of them is equal to 1. The direction of $\tilde{\mathbf{y}}$ is also optimized to be as close as possible to that of $2\hat{\mathbf{y}}' - \hat{\mathbf{y}}$.

Feature Perturbation. Given a query \mathbf{x} with the predicted label $\hat{y} = \arg \max f(\mathbf{x}; \theta)$ that is not in the sensitive region, i.e., $\|g(f^E(\mathbf{x})) - \mathbf{c}^{\hat{y}}\| > r^{\hat{y}}$, the defender creates a falsified confidence vector $\tilde{\mathbf{y}}$ by iteratively perturbing the query feature.

Let $u = f^E(\mathbf{x})$ denote the query feature in \mathcal{F} . The defender first finds out the most distant cluster center in \mathbf{c} to u by computing

$$\arg \max_y \|u - \mathbf{c}^y\|, \quad (19)$$

where $\mathbf{c}^y \in \mathcal{F}$ is computed via Eq. 6. Next, the defender perturbs u by moving u towards \mathbf{c}^y stepwise until the predicted label of u changes. Let $v = \frac{u - \mathbf{c}^y}{\|u - \mathbf{c}^y\|}$ be the vector of movement, the perturbation process is defined as

$$u \leftarrow u + \epsilon v, \text{ s.t. } \arg \max f_{\text{sm}}(f^C(u)) = \hat{y} \quad (20)$$

where ϵ denotes the step size of the perturbation. Eventually, the defender returns $\tilde{y} = f^C(u)$ to the user. Intuitively, feature perturbation aims to obfuscate the decision boundary of the piracy model by making $f(\mathbf{x}; \theta)$ appear at the edge of the two most irrelevant classes, thereby damaging the decision boundary in a similar way compared to boundary unlearning [35].

V. THEORETICAL ANALYSIS

The proofs of the following theorems are given in Section C in the Supplemental Materials.

A. Feasibility of Gradient Reverse

Theorem 1. (The sufficient condition of gradient reverse) To achieve gradient reverse in Definition 4, it is enough to set $\tilde{y} = 2\hat{y}' - \hat{y}$.

B. Certifiability

The aim of the defender is to control the number of honestly answered queries η . Thus, the essential problem is to estimate the maximum number of honestly answered queries $\hat{\eta}$, given the maximum allowable error ϵ with the error probability δ by the defender. In other words, **if the defender pre-defines ϵ and δ , the defender knows how to set the threshold t and radius r** , because η can be estimated by t and r . To achieve this aim, we would like to theoretically identify the relationship between ϵ , δ , t and r with the help of the probably approximately correct (PAC) learning theory [36].

Settings. Suppose the adversary uses its auxiliary dataset D' containing η samples to query the protectee model f and collects the related softmax outputs as ground-truth labels. Within the sensitive region, all softmax outputs are honestly given by f . Let h denote the piracy model trained on D' with respect to the query results at its finest. We consider the true error $E(h)$ as the probability that h makes a mistake on the sample (\mathbf{x}, \mathbf{y}) from the problem domain \mathcal{D} , $E(h) = \Pr_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}}[h(\mathbf{x}) \neq \mathbf{y}]$, and similarly $E(h_e)$ is the empirical error that describes the mistake made by h , $E(h_e) = \frac{1}{\eta} \sum_{\mathbf{x} \in D'} \mathbf{1}[h(\mathbf{x}) \neq f(\mathbf{x})]$. Let t denote the threshold and r refers to the radius. Intuitively, the defender can estimate t and r if he or she has an expectation on how much the adversary can learn from the queries.

Theorem 2. Given a learning algorithm that learns a piracy model h , let η be the actual number of honestly answered sensitive queries, the piracy model can at most be trained to have the maximum allowable error ϵ and the upperbound of error probability δ , if

$$\eta \leq \hat{\eta} = \frac{1}{2\epsilon^2} \cdot \ln\left(\frac{2}{\delta}\right), \quad (21)$$

where $\hat{\eta}$ is the maximum number of honestly answered sensitive queries. Further, let t be the threshold, and r be the query radius, we will have the relationship as follows.

$$r \geq \sqrt{\frac{2t}{\ln(\frac{2}{\delta})}} \epsilon \bar{d}, \quad (22)$$

where \bar{d} is a constant denoting the average distance of features to their cluster center.

Furthermore, we provide another proof that justifies the effectiveness of QUEEN.

Theorem 3. [30], [37] Any adaptive model extraction attack with an arbitrary recovery function $R(\cdot)$ cannot attain a smaller gap between the recovered predictions $R(\tilde{Y}) \in \mathbb{R}^{M \times N}$ and the original predictions $Y \in \mathbb{R}^{M \times N}$ than the following lower bound:

$$\mathbb{E}[\|R(\tilde{Y}) - Y\|_2^2] \geq \frac{MN}{2\pi e} \exp\left(\frac{2}{MN} h(Y|\tilde{Y})\right), \quad (23)$$

where M and N respectively denote the number of the samples and classes; $h(Y|\tilde{Y})$ is the conditional entropy.

VI. EXPERIMENT

A. Experiment Settings

The details of datasets, model architectures, and the implementation of QUEEN are in Section A in the Supplemental Materials. We introduce the attacks and defenses used in the experiment as follows.

Attacks. A model extraction attack consists of two parts: the query strategy and attack strategy. In other words, different attack strategies can be combined with various query strategies.

Two query strategies are considered in this experiment: JBDA-TR [25] and KnockoffNet [14]. KnockoffNet uses only natural data in the auxiliary dataset, whereas JBDA-TR synthesizes data from a small seed set sampled from the auxiliary dataset using Jacobian-based data augmentation [16]. We set the query budget of KnockoffNet to 50,000 and 1,000 for the size of the seed set of JBDA-TR to ensure that it suffices to allow the adversary to get the best attack accuracy. We utilize the random knockoff-net configuration as described in [14], [30].

The employed attack strategies are listed as follows:

- 1) **Direct Query:** The adversary directly uses the output from the protectee model to train the piracy model.
- 2) **Label-Only:** The top-1 hard label is kept, whereas the other results in the softmax output are ignored in the training process.
- 3) **S4L** [38]: Each query image is differently augmented multiple times such that the query results of the different versions of the original query image are averaged to recover the perturbed outputs.
- 4) **Smoothing** [39]: Each query image is augmented by random affine augmentations. Similar to S4L attack, the query results of different augmented versions of the same image are eventually averaged to recover the perturbed outputs.
- 5) **D-DAE** [2]: Meta-classifiers are used to identify the defense method employed by the defender. A number of shadow models are trained to generate clean outputs. A generative model is then trained to recover the perturbed output to the normal output.
- 6) **D-DAE+** [30]: This is the improved version of D-DAE, because it utilizes the lookup table generated for Partial

TABLE II
EVALUATION OF DEFENSES AGAINST ATTACKS ON CIFAR10.

Query Method	Attack Method	None	RS	MAD	AM	Label-only	Rounding	EMDP	ModelGuard	QUEEN
KnockoffNet	Direct Query	87.42%	85.33%	84.58%	83.17%	83.78%	86.77%	66.15%	74.88%	10.00%
	Label-Only	83.78%	83.78%	83.78%	82.11%	83.78%	83.78%	83.78%	83.78%	81.17%
	S4L	86.17%	82.30%	80.21%	82.12%	84.02%	85.86%	66.76%	70.69%	10.00%
	Smoothing	65.43%	63.41%	61.23%	62.27%	61.01%	65.10%	64.36%	53.24%	10.00%
	D-DAE	87.42%	85.32%	84.36%	78.38%	85.24%	87.45%	71.43%	64.73%	78.21%
	D-DAE+	87.42%	85.91%	86.44%	84.51%	84.55%	87.01%	86.43%	58.17%	50.24%
JBDA-TR	pBayes	87.42%	85.91%	87.24%	86.93%	84.57%	86.99%	85.41%	85.16%	84.24%
	Direct Query	63.51%	67.01%	55.31%	60.86%	63.31%	73.55%	25.92%	37.91%	10.00%
	Label-Only	63.51%	63.51%	63.51%	55.77%	63.51%	63.51%	63.51%	63.51%	61.45%
	D-DAE	74.41%	56.63%	48.51%	57.17%	60.15%	67.65%	62.17%	59.17%	47.10%
	D-DAE+	74.41%	72.48%	68.33%	66.10%	63.44%	73.07%	62.33%	51.85%	40.88%
Protectee Model Accuracy	pBayes	74.41%	71.21%	68.15%	74.11%	65.42%	75.01%	67.93%	65.54%	65.33%
	Max Piracy Model Accuracy	87.42%	85.91%	87.24%	86.93%	85.24%	87.45%	86.43%	85.16%	84.24%
	Max Piracy Model Agreement	88.24%	87.21%	89.01%	88.22%	87.13%	88.67%	88.71%	86.78%	86.54%
Protectee Model Accuracy		92.74%	92.74%	92.74%	90.15%	92.74%	92.74%	92.74%	92.74%	90.01%

Bayes Attack (listed below) as the training data of the generative model.

- 7) **Partial Bayes (pBayes)** [30]: This attack employs independent sampling and the neighborhood sampling to establish the Bayes estimator to recover the perturbed outputs.

Defenses. The following defenses are evaluated in our experiments:

- 1) **None**: No defense method is employed. Clean softmax outputs are directly returned to the adversary.
- 2) **Reverse Sigmoid (RS)** [11]: The clean outputs are perturbed by a reverse sigmoid function to enlarge the cross entropy loss.
- 3) **Maximizing Angular Deviation (MAD)** [24]: The clean outputs are perturbed by maximizing the angular deviation between the gradients calculated using the clean outputs and the perturbed outputs.
- 4) **Adaptive Misinformation (AM)** [12]: The protectee model is trained to generate perturbed outputs when it meets Out-Of-Distribution (OOD) queries.
- 5) **Label-Only**: The hard top-1 label is sent back to the user.
- 6) **Rounding**: The confidence score in the softmax output is pruned to have only one decimal place. For example, $[0.14, 0.46, 0.31, 0.09] \leftarrow [0.1, 0.5, 0.3, 0.1]$.
- 7) **Exponential Mechanism Differential Privacy (EMDP)** [40]: The softmax output is perturbed by adding noises generated via exponential mechanism based differential privacy.
- 8) **ModelGuard** [30]: The ModelGuard here refers to ModelGuard-W in the original paper, because ModelGuard-W is better than ModelGuard-S in most cases. The clean output is perturbed by solving a constrained optimization problem that tries to maximize the cross-entropy loss when the adversary trains the piracy model with the perturbed outputs.
- 9) **QUEEN**: The query's sensitivity is measured and the output is perturbed as described in Section IV-A.

B. Experimental Results

In this section, two groups of experiments are presented. First, we evaluate the effectiveness of QUEEN along with the other defenses against various attacks with fixed hyperparameters. It is observed that QUEEN outperforms the other defenses in most cases in the same utility scope. Second, we test how the hyperparameters of QUEEN trade off the performance of the protectee model and the piracy model.

Effectiveness of QUEEN. The defenses are evaluated with the following hyperparameters. ϵ of MAD, ModelGuard is set to be 1.0 such that the accuracy of the protectee model is not affected, while the perturbation is maximized. For the same reason, we set $\gamma = 0.2$ and $\beta = 0.2$ for RS. τ of AM is set to be 0.3 to ensure that the accuracy of the protectee model does not drop significantly. ϵ of EMDP is set to be 1.0. Similarly, we set $r = 0.005$ and $t = 0.2$ for QUEEN to avoid a drastic decrease in accuracy.

The experimental results of effectiveness conducted on CIFAR10 are presented in Table II. The results of the experiments conducted on CIFAR100, Caltech256 and CUB200 are in Table IX, XI and X in the Supplemental Materials. Each row in the table contains the accuracy of the piracy model derived by launching the attack of this row on the defense of the column. At the bottom of the tables, we list the maximum accuracy and agreement of the piracy model among all attacks and the accuracy of the protectee model guarded by the defenses.

Based on the experimental results, we observe that QUEEN outperforms the other state-of-the-art defenses in most cases. For the attacks that do not drastically modify the perturbed output such as Direct Query, S4L and Smoothing, QUEEN can make the piracy model have random-guessing accuracy. Although the accuracy of the protectee model is decreased, this can be considered as the cost of the defense performance. However, when facing the recovery-based attacks, namely D-DAE, D-DAE+, and pBayes, the accuracy of the piracy model increases significantly. Compared to the other defenses, the perturbation scheme of QUEEN is still effective, but not strong enough to completely defend against such attacks. Because the generative model of the recovery-based attacks can bypass the

perturbation-based defenses to an extent.

Following the most recent machine unlearning evaluation metric proposed by Maini et al. [41], we use the forget quality to measure the effectiveness of the defenses. We generate two groups of predicted labels by feeding the query dataset into the piracy model trained with the perturbed output and the piracy model trained with the original output, which are defined as the defended output, and the undefended output. Different from the truth ratio used in [41], we perform Kolmogorov-Smirnov (KS) test to make a statistical test between the two histograms of the defended output and the undefended output to test whether their difference is significant.

As shown in Figure 4, the x and y axes denote the accuracy of the piracy model and the p -value derived from the KS test. The size of the point represents the number of training epoch. It is observed that QUEEN makes the piracy model have only random-guessing level performance when facing direct query attack. The low p -value of Queen reflects that the difference between the output of the piracy models trained with the perturbed and the original query outputs is significant. Notably, ModelGuard also provides a very low p -value.

C. Ablation Study

Impact of the Hyperparameters. We further test how the selection of threshold t is going to affect the performance of QUEEN. We launch KnockoffNet with D-DAE attack on the protectee model trained on CIFAR-10 using TinyImageNet200 as the auxiliary dataset. The query budget is set to 50,000, and t varies in the range of $\{0.1, 0.2, 0.3, 0.4, 0.5\}$ with $r = 0.005$.

The results are shown in Figure 5, where we test how the varying t specifically affects the recorded ratio, reversed ratio, attack accuracy and defense accuracy. For each t value, we repeat the experiment 5 times. The recorded ratio refers to the number of recorded features over the query budget. Similarly, the reversed ratio is the number of queries that are gradient-reversed over the query budget. We omit the results of $t > 0.5$, because a larger t does not lead to any change compared to $t = 0.5$. With $t = 0$, QUEEN performs gradient reverse on every query that hits the sensitive region, which leads to the theoretically lowest attack and defense accuracy. However, this prevents the normal users from getting the honest answer from the protectee model. Hence, we do not consider the $t = 0$ situation in this experiment. It is observed that as t increases, the reversed ratio decreases, and the recorded ratio increases. This means that more queries in the sensitive region are honestly answered, which results in both the ascending attack and defense accuracy. In conclusion, a larger t leads to a lower attack accuracy at the cost of the inevitable decrease in defense accuracy.

Next, we evaluate the impact of the selection of the radius r . We set $r \in \{0.001, 0.005, 0.01, 0.05, 0.1\}$, while keeping $t = 0.2$. The rest of the experiment settings remain the same as above. As depicted in Figure 6, we observe that as r increases, the recorded ratio drops, because the area of the query increases such that the threshold is more quickly reached. For the same reason, one recorded query thus has more neighbor queries, resulting in a decreasing reversed ratio

TABLE III
IMPACT OF THE NUMBER OF SHADOW MODELS ON THE ATTACK ACCURACY OF THE PIRACY MODELS ON CIFAR10 ATTACKED BY D-DAE+.

Number of Models	1	2	3	5	10
Attack Accuracy	52.09%	52.59%	52.56%	51.27%	50.24%
Attack Agreement	52.57%	53.00%	53.13%	52.04%	50.61%

TABLE IV
RUNTIME TEST.

Task	MNIST	Dataset	CIFAR100
		CIFAR10 Runtime (s)	
Training Feature Extraction	12.37	15.30	16.21
Mapping Network Training	63.89	66.91	75.47
Sensitivity Analysis	1.44	1.39	1.47
Training 10 Shadow Models	82.79	94.18	96.41
Process 1,000 Queries	2.51	4.28	4.61

when r is too large. This explains the inverse tendency in reversed ratio, attack accuracy, and defense accuracy when $r = 0.1$. It means that selecting an appropriate r can provide a better trade-off between the performance of the defense and the model's utility. Based on the results, we notice that when $t = 0.2$ and $r = 0.05$, QUEEN provides the most balanced defense, where the defense accuracy does not significantly drop, and the attack accuracy is effectively reduced.

We then investigate how different numbers of shadow models impact QUEEN by testing numbers in $\{1, 2, 3, 5, 10\}$ on the CIFAR10 dataset. The results are detailed in Table III. It is observed that the number of shadow models does not significantly affect the defense effectiveness. We also evaluate the impact of the number of shadow models on the CIFAR10 and CIFAR100 dataset using D-DAE+ and pBayes attacks, and the results are in Table VI, VII and VIII in the Supplemental Materials.

Runtime. As demonstrated in Table IV, we test the time consumption in terms of the preparation and evaluation of QUEEN. The test is conducted over three datasets, and the results of each dataset are close to each other. This is because of the number of samples in these training datasets are equal to 60,000. It is observed that the training feature extraction takes no longer than 17 seconds, and training the mapping network costs less than 76 seconds. Additionally, the time of sensitivity analysis and query processing is trivial. The comparison of time cost between the defenses are listed in Table XII in the Supplemental Materials. Overall, the runtime is practical for QUEEN to be applied in real-world scenario.

VII. DISCUSSION

The effectiveness of QUEEN is established upon the two main components: sensitivity measurement and output perturbation. By launching gradient reverse on the sensitive query and feature perturbation on the non-sensitive query, QUEEN achieves the goal where the attack accuracy of the piracy model is significantly lowered while the defense accuracy remaining high. However, this defense essentially modifies the

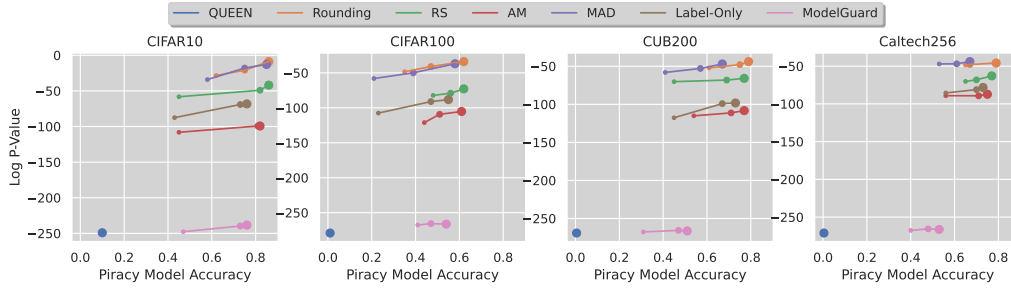
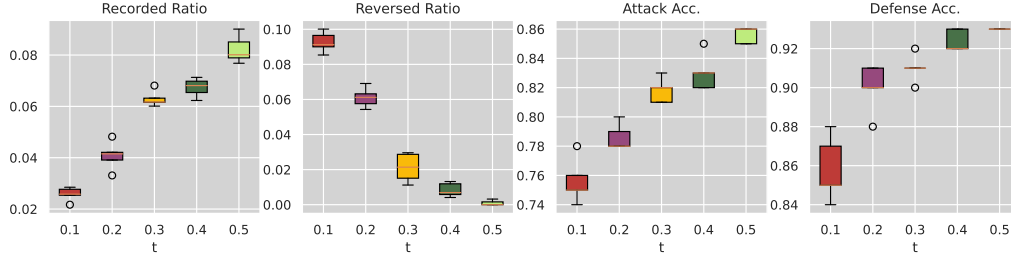
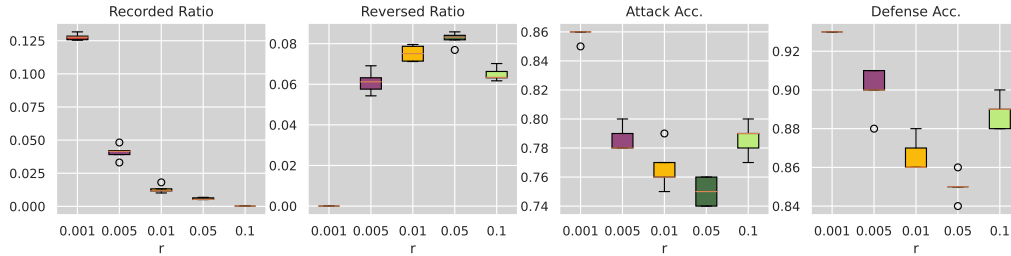


Fig. 4. The results of KS test on the outputs of the piracy models.

Fig. 5. The impact of selection of the threshold t on CIFAR-10.Fig. 6. The impact of selection of the radius r on CIFAR-10.

posterior outputs of the protectee model. In fact, after receiving the outputs from the defender, the adversary is allowed to take any operation on the paired auxiliary dataset. This means that if the adversary uses advanced attacks such as D-DAE to purify the perturbed outputs with the generator network, the attack accuracy will be lifted. QUEEN is capable of defending such advanced attacks to an extent, but the effectiveness is limited if the generator used for purifying the outputs can make the argmax of the softmax output correct. To mitigate the advanced attacks could be a future research direction. However, to launch the advanced attacks, the adversary is required to train a large amount of shadow models, meta-classifiers, whose cost is non-negligible.

Currently, the mainstream defenses [11], [24], [12], [30] counter the model extraction attacks based on perturbation of the softmax outputs. In these studies, the adversary is assumed to train the piracy model with the softmax outputs, because this leads to better test performance. Using multiple accounts to decrease the defense strength can be countered by IP detection defenses. A more complex attack is feasible such as querying the model in the distributed denial-of-service (DDOS) manner, but the cost of the attack will significantly increase.

VIII. CONCLUSION

We propose QUEEN, a proactive defense against model extraction attacks by detecting potential threats of the queries by measuring the cumulative query sensitivity. The softmax of the query whose feature is within the sensitive region is perturbed to make gradient reverse if the cumulative query sensitivity exceeds the threshold. The features of the non-sensitive queries are perturbed so as to generate perturbed softmax. Through extensive experiments, the effectiveness of QUEEN has been proved, where QUEEN is capable of defending the current model extraction attacks and has outperformed the SOTA defenses. The attack accuracy is decreased at an acceptable cost of the defense accuracy.

REFERENCES

- [1] R. Gozalo-Brizuela and E. C. Garrido-Merchan, "Chatgpt is not all you need. a state of the art review of large generative ai models," *arXiv preprint arXiv:2301.04655*, 2023.
- [2] Y. Chen, R. Guan, X. Gong, J. Dong, and M. Xue, "D-dae: Defense-penetrating model extraction attacks," in *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2023, pp. 382–399.
- [3] C. A. Choquette-Choo, F. Tramèr, N. Carlini, and N. Papernot, "Label-only membership inference attacks," in *Proceedings of the 38th International Conference on Machine Learning*, 2021, pp. 1964–1974.

- [4] Y. Liu, Z. Zhao, M. Backes, and Y. Zhang, "Membership inference attacks by exploiting loss trajectory," in *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, 2022, pp. 2085–2098.
- [5] T. Cong, X. He, and Y. Zhang, "Sslguard: A watermarking scheme for self-supervised learning pre-trained encoders," in *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, 2022, pp. 579–593.
- [6] H. Jia, C. A. Choquette-Choo, V. Chandrasekaran, and N. Papernot, "Entangled watermarks as a defense against model extraction," in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 1937–1954.
- [7] N. Lukas, Y. Zhang, and F. Kerschbaum, "Deep neural network fingerprinting by conferrable adversarial examples," *arXiv preprint arXiv:1912.00888*, 2019.
- [8] Z. Peng, S. Li, G. Chen, C. Zhang, H. Zhu, and M. Xue, "Fingerprinting deep neural networks globally via universal adversarial perturbations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 13 430–13 439.
- [9] J. Guan, J. Liang, and R. He, "Are you stealing my model? sample correlation for fingerprinting deep neural networks," *Advances in Neural Information Processing Systems*, vol. 35, pp. 36 571–36 584, 2022.
- [10] A. Dziedzic, M. A. Kaleem, Y. S. Lu, and N. Papernot, "Increasing the cost of model extraction with calibrated proof of work," in *International Conference on Learning Representations*, 2022.
- [11] T. Lee, B. Edwards, I. Molloy, and D. Su, "Defending against machine learning model stealing attacks using deceptive perturbations," *arXiv preprint arXiv:1806.00054*, 2018.
- [12] S. Kariyappa and M. K. Qureshi, "Defending against model stealing attacks with adaptive misinformation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 770–778.
- [13] J. R. Correia-Silva, R. F. Berriel, C. Badue, A. F. de Souza, and T. Oliveira-Santos, "Copicat cnn: Stealing knowledge by persuading confession with random non-labeled data," in *2018 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2018, pp. 1–8.
- [14] T. Orekondy, B. Schiele, and M. Fritz, "Knockoff nets: Stealing functionality of black-box models," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 4954–4963.
- [15] S. Pal, Y. Gupta, A. Shukla, A. Kanade, S. Shevade, and V. Ganapathy, "Activethief: Model extraction using active learning and unannotated public data," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 01, 2020, pp. 865–872.
- [16] N. Papernot, P. McDaniel, I. Goodfellow, S. Jha, Z. B. Celik, and A. Swami, "Practical black-box attacks against machine learning," in *Proceedings of the 2017 ACM on Asia conference on computer and communications security*, 2017, pp. 506–519.
- [17] H. Yu, K. Yang, T. Zhang, Y.-Y. Tsai, T.-Y. Ho, and Y. Jin, "Cloudleak: Large-scale deep learning models stealing through adversarial examples," in *NDSS*, 2020.
- [18] J.-B. Truong, P. Maini, R. J. Walls, and N. Papernot, "Data-free model extraction," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 4771–4780.
- [19] S. Kariyappa, A. Prakash, and M. K. Qureshi, "Maze: Data-free model stealing attack using zeroth-order gradient estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 13 814–13 823.
- [20] Y. Li, M. Zhu, X. Yang, Y. Jiang, T. Wei, and S.-T. Xia, "Black-box dataset ownership verification via backdoor watermarking," *IEEE Transactions on Information Forensics and Security*, 2023.
- [21] H. Zhang, G. Hua, X. Wang, H. Jiang, and W. Yang, "Categorical inference poisoning: Verifiable defense against black-box dnn model stealing without constraining surrogate data and query times," *IEEE Transactions on Information Forensics and Security*, vol. 18, pp. 1473–1486, 2023.
- [22] W. Peng, J. Yi, F. Wu, S. Wu, B. Zhu, L. Lyu, B. Jiao, T. Xu, G. Sun, and X. Xie, "Are you copying my model? protecting the copyright of large language models for eaas via backdoor watermark," *arXiv preprint arXiv:2305.10036*, 2023.
- [23] A. Hu, Z. Lu, R. Xie, and M. Xue, "Veridip: Verifying ownership of deep neural networks through privacy leakage fingerprints," *IEEE Transactions on Dependable and Secure Computing*, 2023.
- [24] T. Orekondy, B. Schiele, and M. Fritz, "Prediction poisoning: Towards defenses against dnn model stealing attacks," in *8th International Conference on Learning Representations*, 2020.
- [25] M. Juuti, S. Szyller, S. Marchal, and N. Asokan, "Prada: protecting against dnn model stealing attacks," in *2019 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2019, pp. 512–527.
- [26] S. Kariyappa, A. Prakash, and M. K. Qureshi, "Protecting dnns from theft using an ensemble of diverse models," in *International Conference on Learning Representations*, 2020.
- [27] J. Zhang, S. Peng, Y. Gao, Z. Zhang, and Q. Hong, "Apmsa: adversarial perturbation against model stealing attacks," *IEEE Transactions on Information Forensics and Security*, vol. 18, pp. 1667–1679, 2023.
- [28] D. Hendrycks, S. Basart, M. Mazeika, A. Zou, J. Kwon, M. Mostajabi, J. Steinhardt, and D. Song, "Scaling out-of-distribution detection for real-world settings," in *International Conference on Machine Learning*. PMLR, 2022, pp. 8759–8773.
- [29] J. Ye, A. Maddi, S. K. Murakonda, V. Bindschaedler, and R. Shokri, "Enhanced membership inference attacks against machine learning models," in *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, 2022, pp. 3093–3106.
- [30] M. Tang, A. Dai, L. DiValentin, A. Ding, A. Hass, N. Z. Gong, and Y. Chen, "Modelguard: Information-theoretic defense against model extraction attacks."
- [31] L. Deng, "The mnist database of handwritten digit images for machine learning research [best of the web]," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- [32] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, "Supervised contrastive learning," *Advances in neural information processing systems*, vol. 33, pp. 18 661–18 673, 2020.
- [33] S. Zhang, R. Xu, C. Xiong, and C. Ramaiah, "Use all the labels: A hierarchical multi-label contrastive learning framework," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16 660–16 669.
- [34] N. Haim, G. Vardi, G. Yehudai, O. Shamir, and M. Irani, "Reconstructing training data from trained neural networks," *Advances in Neural Information Processing Systems*, vol. 35, pp. 22 911–22 924, 2022.
- [35] M. Chen, W. Gao, G. Liu, K. Peng, and C. Wang, "Boundary unlearning: Rapid forgetting of deep networks via shifting the decision boundary," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7766–7775.
- [36] L. G. Valiant, "A theory of the learnable," *Communications of the ACM*, vol. 27, no. 11, pp. 1134–1142, 1984.
- [37] T. M. Cover, *Elements of information theory*. John Wiley & Sons, 1999.
- [38] M. Jagielski, N. Carlini, D. Berthelot, A. Kurakin, and N. Papernot, "High accuracy and high fidelity extraction of neural networks," in *29th USENIX security symposium (USENIX Security 20)*, 2020, pp. 1345–1362.
- [39] N. Lukas, E. Jiang, X. Li, and F. Kerschbaum, "Sok: How robust is image classification deep neural network watermarking?" in *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2022, pp. 787–804.
- [40] C. Ilvento, "Implementing the exponential mechanism with base-2 differential privacy," in *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, 2020, pp. 717–742.
- [41] P. Maini, Z. Feng, A. Schwarzschild, Z. C. Lipton, and J. Z. Kolter, "Tofu: A task of fictitious unlearning for llms," *arXiv preprint arXiv:2401.06121*, 2024.
- [42] G. Griffin, A. Holub, and P. Perona, "Caltech 256," April 2022.
- [43] C. Wah, S. Branson, P. Welinder, P. Perona, and S. J. Belongie, "The caltech-ucsd birds-200-2011 dataset," 2011. [Online]. Available: <https://api.semanticscholar.org/CorpusID:16119123>
- [44] A. Krizhevsky, "Learning multiple layers of features from tiny images," 2009. [Online]. Available: <https://api.semanticscholar.org/CorpusID:18268744>
- [45] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [46] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [47] M. A. mnmostafa, "Tiny imagenet," 2017. [Online]. Available: <https://kaggle.com/competitions/tiny-imagenet>
- [48] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.

TABLE V
EXPERIMENT SETTINGS

Defender	Training Dataset	CIFAR10	CIFAR100	CUB200	Caltech256
	OE Dataset	SVHN	SVHN	Indoor67	Indoor67
	Model Architecture	VGG16-BN	VGG16-BN	ResNet50	ResNet50
	Accuracy	92.74%	76.12%	82.42%	86.84%
Adversary	Auxiliary Dataset	TinyImageNet200	TinyImageNet200	ImageNet1k	ImageNet1k
	Model Architecture	VGG16-BN	VGG16-BN	ResNet50	ResNet50

APPENDIX

A. Experiment Settings

Datasets and Model Architectures. The information of the training datasets and the corresponding model architectures is summarized in Table V. Four types of model architectures are used for training the protectee models on four image classification datasets, which are Caltech256 [42], CUB200[43], CIFAR100, and CIFAR10 [44]. For the defender, the protectee models are all trained with Outlier Exposure (OE) datasets, because one of the defense methods, Adaptive Misinformation [12], requires it. For Caltech256 and CUB200, the defender employs ResNet50 [45]. For CIFAR100 and CIFAR10, the defender uses VGG16-BN [46]. The test accuracy of the protectee models is presented in Table V. For the adversary, the same model architectures are selected to allow the adversary to obtain the best attack accuracy. In terms of the auxiliary datasets, the adversary use TinyImageNet200 [47] for CIFAR10 and CIFAR100, and ImageNet1k [48] for CUB200 and Caltech256.

Evaluation Metric. In the experiment, we mainly use attack accuracy and attack agreement as the evaluation metrics for the attack performance. Attack accuracy is defined as the ratio of the correctly classified samples to the total samples in the test dataset of the protectee model. Attack agreement, on the other hand, is the ratio of the samples that are identically classified by both the protectee and piracy models to the total samples in the test dataset. The lower the values of the two metrics are, the better the defense is.

Implementation of QUEEN For the mapping network, we use four fully connected layers to map the training features to the 2D space. As for the piracy model ensemble, we select ResNet18 for Caltech256 and CUB200 datasets, and VGG11-BN for the CIFAR-10 and CIFAR-100 datasets. The mapping network is trained for 100 epochs using the SGD optimizer, where the learning rate is set to 0.01 and it decreases by 0.5 every 20 epochs. The members of the piracy model ensemble are trained for 5 epochs on the sub-datasets in each of which there are 500 samples in each class.

B. Peripheral Data vs. Central Data.

We split the training dataset of MNIST to train models and test their accuracy. The cluster center of each class is computed such that the samples in each class is ranked based on the distance between their features and the cluster centers. The closer the query feature to the cluster center is, the lower the rank of the feature will be. We split the training dataset into

TABLE VI
IMPACT OF THE NUMBER OF SHADOW MODELS ON THE ATTACK ACCURACY OF THE PIRACY MODELS ON CIFAR10 ATTACKED BY PBAYES.

Number of Models	1	2	3	5	10
Attack Accuracy	67.12%	66.41%	65.84%	65.33%	65.67%
Attack Agreement	69.41%	68.57%	67.76%	67.56%	67.81%

TABLE VII
IMPACT OF THE NUMBER OF SHADOW MODELS ON THE ATTACK ACCURACY OF THE PIRACY MODELS ON CIFAR100 ATTACKED BY D-DAE+.

Number of Models	1	2	3	5	10
Attack Accuracy	57.57%	56.78%	57.12%	56.59%	56.61%
Attack Agreement	59.46%	58.89%	59.25%	59.11%	58.57%

TABLE VIII
IMPACT OF THE NUMBER OF SHADOW MODELS ON THE ATTACK ACCURACY OF THE PIRACY MODELS ON CIFAR100 ATTACKED BY PBAYES.

Number of Models	1	2	3	5	10
Attack Accuracy	57.32%	56.51%	56.72%	56.84%	56.45%
Attack Agreement	59.81%	59.04%	59.31%	58.95%	58.67%

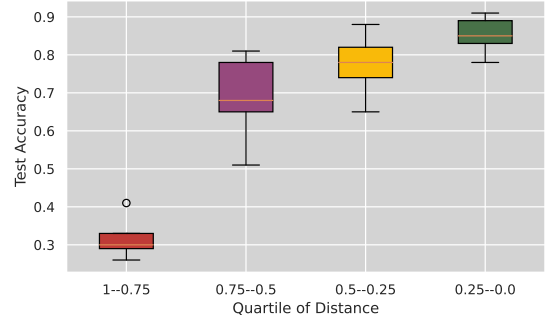


Fig. 7. Accuracy of Models Trained on Peripheral and Central Data.

four quartiles based on the rank. For example, the 0.25–0.0 quartile is the subset containing the top 25% of the training samples that are closest to the cluster centers. In each training attempt, 500 samples are randomly sampled from each class of the quartile, and then used for training the model. For each quartile, we repeat for five times to get the results, where the settings remain the same as in training the protectee models on MNIST in the previous experiments.

As depicted in Figure 7, the results support our assumption,

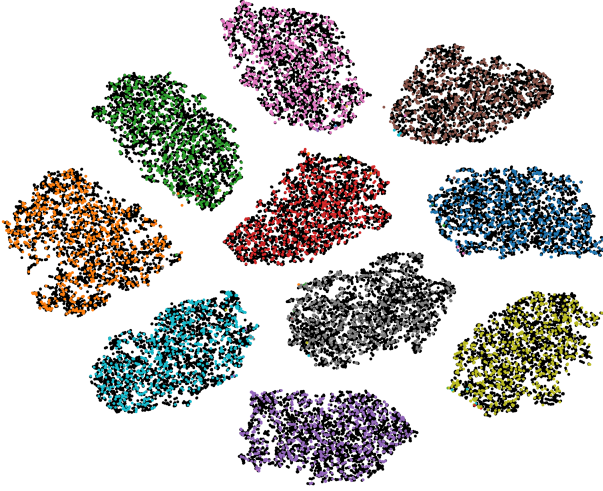


Fig. 8. Visualization of the decision boundary via t-SNE. The black/colorful dots denote the features of the training/test data, where each color represents one class.

where the peripheral subset leads to significantly lower test accuracy than the central subset. This suggests that the central data samples are more valuable than the peripheral data samples in terms of protection. Output perturbation that will lower the defense accuracy such as gradient reverse should be applied to the central data, because the attack accuracy can be efficiently lowered at a less cost of defense accuracy. This explains the idea of QUEEN: punishing the critical queries and letting the trivial queries pass.

C. Supplementary of Theoretical Analysis

1) *Proof of Theorem 1: **Theorem 1.** (The sufficient condition of gradient reverse) To achieve gradient reverse in Definition 4, it is enough to set $\tilde{\mathbf{y}} = 2\hat{\mathbf{y}}' - \hat{\mathbf{y}}$.*

Proof. We consider the widely used cross-entropy loss here $\mathcal{L}_{CE} = -\sum_{i=1}^n \mathbf{y}^i \log(\hat{\mathbf{y}}^i)$, where \mathbf{y}^i is the i -th element of the ground-truth label vector \mathbf{y} .

According to the chain rule,

$$\nabla_{\tau} \mathcal{L}_{CE} = \frac{\partial \mathcal{L}_{CE}}{\partial \hat{\mathbf{y}}'} \cdot \frac{\partial \hat{\mathbf{y}}'}{\partial \mathbf{h}} \cdot \frac{\partial \mathbf{h}}{\partial \tau}, \quad (24)$$

where $\hat{\mathbf{y}}' = f_{sm}(h(\mathbf{x}'))$.

The gradient reverse in Definition 4 requires that

$$\frac{\partial \mathcal{L}_{CE}}{\partial \tilde{\mathbf{y}}} \cdot \frac{\partial \tilde{\mathbf{y}}}{\partial \mathbf{h}} \cdot \frac{\partial \mathbf{h}}{\partial \tau} = -\frac{\partial \mathcal{L}_{CE}}{\partial \hat{\mathbf{y}}'} \cdot \frac{\partial \hat{\mathbf{y}}'}{\partial \mathbf{h}} \cdot \frac{\partial \mathbf{h}}{\partial \tau}. \quad (25)$$

Eventually, we have

$$\tilde{\mathbf{y}} = 2\hat{\mathbf{y}}' - \hat{\mathbf{y}}. \quad (26)$$

That completes the proof. \square

2) *Proof of Theorem 2: **Theorem 2.** Given a learning algorithm that learns a piracy model h , let η be the actual number of honestly answered sensitive queries, the piracy model can at most be trained to have the maximum allowable error ϵ and the upperbound of error probability δ , if*

$$\eta \leq \hat{\eta} = \frac{1}{2\epsilon^2} \cdot \ln\left(\frac{2}{\delta}\right), \quad (27)$$

where $\hat{\eta}$ is the maximum number of honestly answered sensitive queries. Further, let t be the threshold, and r be the query radius, we will have the relationship as follows.

$$r \geq \sqrt{\frac{2t}{\ln(\frac{2}{\delta})}} \epsilon \bar{d}, \quad (28)$$

where \bar{d} is a constant denoting the average distance of features to their cluster center.

Proof. According to hoeffding's inequality, we have

$$Pr\left(|E(h) - E(h_e)| > \epsilon\right) \leq 2e^{-2\hat{\eta}\epsilon^2}. \quad (29)$$

Let $\delta = 2e^{-2\hat{\eta}\epsilon^2}$, we then have

$$\hat{\eta} = \frac{1}{2\epsilon^2} \cdot \ln\left(\frac{2}{\delta}\right). \quad (30)$$

Thus, the piracy model can at most achieve ϵ error rate with δ probability if $\eta \leq \hat{\eta}$.

So far, we have determined the maximum number of honestly answered queries η by ϵ and δ .

In addition, η can be estimated by t and r :

$$\eta = \frac{t\pi\bar{d}^2}{\pi r^2} = \frac{t\bar{d}^2}{r^2} \quad (31)$$

Combine equations 30 and 31 with the condition $\eta \leq \hat{\eta}$, we have

$$\frac{t\bar{d}^2}{r^2} \leq \frac{1}{2\epsilon^2} \cdot \ln\left(\frac{2}{\delta}\right) \quad (32)$$

Eventually, we have

$$r \geq \sqrt{\frac{2t}{\ln(\frac{2}{\delta})}} \epsilon \bar{d} \quad (33)$$

This completes the proof. \square

3) *Proof of Theorem 3: **Theorem 3.** [30], [37] Any adaptive model extraction attack with an arbitrary recovery function $R(\cdot)$ cannot attain a smaller gap between the recovered predictions $R(\tilde{Y}) \in \mathbb{R}^{M \times N}$ and the original predictions $Y \in \mathbb{R}^{M \times N}$ than the following lower bound:*

$$\mathbb{E}[\|R(\tilde{Y}) - Y\|_2^2] \geq \frac{MN}{2\pi e} \exp\left(\frac{2}{MN} h(Y|\tilde{Y})\right), \quad (34)$$

where M and N respectively denote the number of the samples and classes; $h(Y|\tilde{Y})$ is the conditional entropy.

Proof. The following inequality holds for \mathbf{y} with an arbitrary distribution conditioned on the event $\{\tilde{Y} = \tilde{Y}_p\}$ [37]:

$$\begin{aligned} h(Y|\tilde{Y} = \tilde{Y}_p) &\leq \frac{1}{2} \log((2\pi e)^{MN} \det(\text{Cov}(Y|\tilde{Y} = \tilde{Y}_p))) \\ \Rightarrow \det(\text{Cov}(Y|\tilde{Y} = \tilde{Y}_p)) &\geq \frac{1}{2\pi e^{MN}} \exp(2h(Y|\tilde{Y} = \tilde{Y}_p)) \end{aligned} \quad (35)$$

where the equality holds with Gaussian $Y|\{\tilde{Y} = \tilde{Y}_p\}$; \tilde{Y}_p denotes the perturbed outputs generated by the perturbation function given the training labels. With the fact that any

TABLE IX
EVALUATION OF DEFENSES AGAINST ATTACKS ON CIFAR100.

Query Method	Attack Method	None	RS	MAD	AM	Label-only	Rounding	EMDP	ModelGuard	QUEEN
KnockoffNet	Direct Query	66.74%	62.55%	58.79%	60.75%	54.47%	62.91%	48.14%	51.98%	1.0%
	Label-Only	54.47%	54.47%	54.47%	52.61%	54.47%	54.47%	54.47%	54.47%	53.47%
	S4L	63.21%	57.14%	55.39%	60.13%	55.18%	61.75%	50.18%	46.42%	1.0%
	Smoothing	64.27%	63.41%	60.17%	62.53%	61.43%	66.18%	55.81%	51.84%	1.0%
	D-DAE	66.74%	63.19%	62.51%	61.34%	58.61%	63.31%	62.46%	59.10%	58.37%
	D-DAE+	66.74%	64.77%	64.37%	62.10%	58.91%	63.17%	62.66%	57.01%	56.61%
JBDA-TR	pBayes	66.74%	65.55%	65.01%	64.97%	58.18%	65.47%	58.41%	59.11%	56.45%
	Direct Query	41.41%	35.17%	13.35%	30.61%	22.32%	36.11%	9.43%	5.87%	1.0%
	Label-Only	22.32%	22.32%	22.32%	15.32%	22.32%	22.32%	22.32%	22.32%	20.36%
	D-DAE	41.41%	30.18%	20.54%	26.63%	19.78%	24.71%	24.76%	19.31%	18.36%
	D-DAE+	41.41%	38.87%	22.24%	31.44%	21.68%	39.96%	29.14%	20.62%	17.18%
	pBayes	41.41%	39.52%	27.46%	40.04%	24.11%	39.75%	26.91%	22.61%	20.75%
Max Piracy Model Accuracy		66.74%	65.55%	65.01%	64.97%	61.43%	66.18%	62.66%	59.11%	58.37%
Max Piracy Model Agreement		72.57%	71.15%	71.24%	71.48%	66.91%	68.41%	64.28%	65.82%	65.51%
Protectee Model Accuracy		76.12%	76.12%	76.12%	74.45%	76.12%	76.12%	76.12%	76.12%	74.23%

TABLE X
EVALUATION OF DEFENSES AGAINST ATTACKS ON CUB200.

Query Method	Attack Method	None	RS	MAD	AM	Label-only	Rounding	EMDP	ModelGuard	QUEEN
KnockoffNet	Direct Query	80.79%	75.41%	67.15%	77.14%	73.11%	79.91%	71.47%	54.02%	0.5%
	Label-Only	73.11%	73.11%	73.11%	68.15%	73.11%	73.11%	73.11%	73.11%	69.88%
	S4L	80.14%	74.76%	60.55%	75.84%	76.51%	78.47%	72.78%	54.15%	0.5%
	Smoothing	79.91%	75.04%	67.11%	74.52%	75.51%	79.35%	75.88%	51.35%	0.5%
	D-DAE	80.79%	78.35%	78.76%	76.49%	72.18%	79.81%	76.98%	65.47%	70.19%
	D-DAE+	80.79%	79.31%	79.89%	77.40%	71.83%	79.61%	77.21%	74.26%	68.34%
JBDA-TR	pBayes	80.79%	80.28%	79.38%	78.81%	70.58%	80.31%	77.57%	76.47%	69.14%
	Direct Query	64.23%	54.15%	10.14%	36.17%	30.66%	50.91%	14.22%	5.14%	0.5%
	Label-Only	30.66%	30.66%	30.66%	23.70%	30.66%	30.66%	30.66%	30.66%	23.34%
	D-DAE	64.23%	53.25%	16.73%	38.19%	25.72%	44.29%	32.91%	8.91%	15.81%
	D-DAE+	64.23%	62.78%	34.79%	41.74%	33.18%	48.85%	36.48%	28.37%	17.21%
	pBayes	64.23%	62.39%	34.74%	61.69%	32.43%	50.03%	33.17%	27.17%	18.55%
Max Piracy Model Accuracy		80.79%	80.28%	79.89%	78.81%	76.51%	80.31%	77.57%	76.47%	70.19%
Max Piracy Model Agreement		84.21%	84.02%	83.57%	80.28%	79.04%	81.35%	78.94%	78.21%	73.47%
Protectee Model Accuracy		82.42%	82.42%	82.42%	80.19%	82.42%	82.42%	82.42%	82.42%	80.71%

recovery function $R(\cdot)$ tries to minimize $\mathbb{E}[\|R(\tilde{Y}) - Y\|_2^2 | \tilde{Y}]$,
we then have:

$$\begin{aligned}
& \mathbb{E}[\|R(\tilde{Y}) - Y\|_2^2 | \tilde{Y} = \tilde{Y}_p] \\
& \geq \mathbb{E}[\|Y - \mathbb{E}[Y | \tilde{Y} = \tilde{Y}_p]\|_2^2 | \tilde{Y} = \tilde{Y}_p] \\
& = \text{tr}(\text{Cov}(Y | \tilde{Y} = \tilde{Y}_p)) \\
& \geq NC[\det(\text{Cov}(Y | \tilde{Y} = \tilde{Y}_p))]^{\frac{1}{MN}} \\
& \geq \frac{MN}{2\pi e} \exp\left(\frac{2}{MN} h(Y | \tilde{Y} = \tilde{Y}_p)\right).
\end{aligned} \tag{36}$$

That completes the proof. \square

TABLE XI
EVALUATION OF DEFENSES AGAINST ATTACKS ON CALTECH256.

Query Method	Attack Method	None	RS	MAD	AM	Label-only	Rounding	EMDP	ModelGuard	QUEEN
KnockoffNet	Direct Query	83.57%	77.24%	67.21%	75.63%	72.98%	78.81%	69.74%	53.37%	0.39%
	Label-Only	72.98%	72.98%	72.98%	68.33%	72.98%	72.98%	72.98%	72.98%	69.41%
	S4L	80.43%	75.86%	62.01%	75.82%	72.17%	77.49%	68.41%	54.73%	0.39%
	Smoothing	81.77%	75.32%	69.14%	73.55%	76.04%	77.41%	75.93%	55.76%	0.39%
	D-DAE	83.57%	80.25%	79.76%	78.35%	72.68%	75.04%	73.92%	68.31%	66.26%
	D-DAE+	83.57%	81.17%	81.43%	78.65%	70.33%	79.17%	73.21%	70.65%	66.45%
	pBayes	83.57%	81.75%	82.21%	82.65%	72.65%	82.43%	76.14%	78.34%	72.44%
JBDA-TR	Direct Query	64.69%	54.51%	9.17%	38.92%	31.25%	44.21%	12.54%	5.18%	0.39%
	Label-Only	28.91%	28.91%	28.91%	24.91%	28.91%	28.91%	28.91%	28.91%	23.15%
	D-DAE	64.69%	55.57%	25.79%	37.21%	29.31%	45.19%	48.17%	15.26%	29.76%
	D-DAE+	64.69%	62.41%	38.79%	44.99%	38.08%	55.13%	49.56%	37.76%	37.45%
	pBayes	64.69%	61.21%	38.77%	59.74%	36.71%	58.62%	41.73%	35.13%	32.51%
Max Piracy Model Accuracy		83.57%	81.75%	82.21%	82.65%	76.04%	82.43%	76.14%	78.34%	72.44%
Max Piracy Model Agreement		86.14%	85.89%	86.01%	85.13%	79.14%	85.22%	79.34%	83.02%	75.63%
Protectee Model Accuracy		86.84%	86.84%	86.84%	84.24%	86.84%	86.84%	86.84%	86.84%	83.16%

TABLE XII
TIME OF PROCESSING 1,000 QUERIES USING EACH DEFENSE ON CIFAR10.

Defenses	None	RS	MAD	AM	Label-only	Rounding	EMDP	ModelGuard	QUEEN
Time (Seconds)	1.28	1.56	30.76	2.37	1.21	1.37	1.52	2.47	4.28