# Towards Undistillable Models by Minimizing Conditional Mutual Information

Linfeng Ye\*

Department of Electrical and Computer Engineering University of Waterloo

Shayan Mohajer Hamidi\* Department of Electrical and Computer Engineering University of Waterloo

**En-Hui Yang** Department of Electrical and Computer Engineering University of Waterloo

Reviewed on OpenReview: https://openreview.net/forum?id=jVABSsD4Vf

l44ye@uwaterloo.ca

smohajer@uwaterloo.ca

ehy ang @uwaterloo.ca

### Abstract

A deep neural network (DNN) is said to be undistillable if, when used as a black-box inputoutput teacher, it cannot be distilled through knowledge distillation (KD). In this case, the distilled student (referred to as the knockoff student) does not outperform a student trained independently with label smoothing (LS student) in terms of prediction accuracy. To protect intellectual property of DNNs, it is desirable to build undistillable DNNs. To this end, it is first observed that an undistillable DNN may have the trait that each cluster of its output probability distributions in response to all sample instances with the same label should be highly concentrated to the extent that each cluster corresponding to each label should ideally collapse into one probability distribution. Based on this observation and by measuring the concentration of each cluster in terms of conditional mutual information (CMI), a new training method called CMI minimized (CMIM) method is proposed, which trains a DNN by jointly minimizing the conventional cross entropy (CE) loss and the CMI values of all temperature scaled clusters across the entire temperature spectrum. The resulting CMIM model is shown, by extensive experiments, to be undistillable by all tested KD methods existing in the literature. That is, the knockoff students distilled by these KD methods from the CMIM model underperform the respective LS students. In addition, the CMIM model is also shown to perform better than the model trained with the CE loss alone in terms of their own prediction accuracy. The code for the paper is publicly available at https://anonymous.4open.science/r/CMIM-605C/README.md.

### 1 Introduction

Originally aiming for model compression, knowledge distillation (Buciluǎ et al., 2006; Hinton et al., 2015) (KD) has received significant attention from both academia and industry in recent years due to its remarkable effectiveness. The essence of KD is to transfer the knowledge of a pre-trained large model (teacher) to a smaller model (student). Building on the work of Hinton et al. (2015), numerous follow-up works have endeavored to enhance the performance of KD (Romero et al., 2014; Anil et al., 2018; Park et al., 2019) and to gain deeper insights into why distillation is effective (Phuong & Lampert, 2019; Mobahi et al., 2020; Ye et al., 2024; Allen-Zhu & Li, 2020; Menon et al., 2021; Borup & Andersen, 2021).

<sup>\*</sup>Authors contributed equally.

In the scenario where the teacher does not want its knowledge to be transferred, however, KD is undesirable and indeed poses a threat to intellectual property (IP) of the teacher (Shokri & Shmatikov, 2015). Developing and training a high-quality large DNN requires significant investments of time, effort, finances, and resources, including extensive data annotation and computational infrastructure. Developers of the large DNN may want to prevent the knowledge of the large DNN from being transferred by their competitors. However, once the large DNN is released as a "black box", anyone can apply a logit-based KD method (or equivalently a distribution-based KD method (Zheng & YANG, 2024)) to distill the DNN as a teacher. The goal is to train a student, referred to as a knockoff student in the context of DNN IP protection, that mimics the teacher's behavior to gain competitive advantages. As such, in this case, it would be desirable for the developers to build the large DNN so that it is undistillable. The question, of course, is how.

Before delving deeper into the above question, let us first clarify what we mean by saying that a DNN is undistillable. At this point, we invoke the concept of distillable DNN introduced recently in Yang & Ye (2024):

**Definition 1.** [Distillability of a DNN (Yang & Ye, 2024)] When used as a black-box input-output teacher, a DNN is said to be distillable<sup>1</sup> with respect to a student if there exists a KD method which, when applied to the teacher and student, yields a knockoff student outperforming the student trained alone with label smoothing (LS student) in terms of prediction accuracy.

Therefore, a DNN is undistillable if no knockoff student can outperform the respective LS student regardless of which logit-based KD method is used. Since there are so many logit-based KD methods and so many students, what type of a DNN is undistillable? In comparison with the training of LS student, Definition 1 provides some insight into a trait that an undistillable DNN may possess. For each label, consider the cluster of the output probability distributions of the DNN in response to all input sample instances with that label. If each cluster corresponding to each label is highly concentrated to the extent that all probability distributions within the cluster more or less collapse into one probability distribution, then the student training within KD is similar to that of the respective LS student regardless of which logit-based KD method and which student are applied. In this case, one would expect that no knockoff student would perform significantly better than the respective LS student. Therefore, a DNN possessing this trait will likely be undistillable.

From a theoretical perspective, KD relies on the diversity within the teacher's predicted distributions to convey richer contextual information to the student. The greater the diversity in the cluster of predictions for a given class, the more nuanced guidance the student can receive (see Ye et al. (2024)). Conversely, when the teacher's outputs become highly concentrated—meaning predictions for inputs with the same label converge to a single probability vector—this contextual richness is lost. The student, in this case, receives largely redundant information, rendering the distillation process as ineffective as direct training with label smoothing.

To quantify this notion of concentration, we turn to conditional mutual information (CMI) (Yang et al., 2025). Specifically, let X denote the random input sample to the DNN, and Y be the ground truth label of X. Let  $\hat{Y}$  denote the random label predicted by the DNN in response to input X. It was shown in Yang et al. (2025) that for each label y, the label specific CMI  $I(X; \hat{Y} | Y = y)$  measures the concentration of the cluster corresponding to label y, and the CMI  $I(X; \hat{Y} | Y)$  measures the average concentration across all clusters. To build an undistillable DNN, one then is motivated to minimize jointly the conventional cross entropy (CE) loss and the CMI  $I(X; \hat{Y} | Y)$ .

In this paper, we will go one step further. In KD (Hinton et al., 2015), temperature scaling of logits is often applied. It was shown in Zheng & YANG (2024) that logit temperature scaling with temperature T can be equivalently achieved by power transform of the output probability distribution with power  $\alpha = 1/T$ . Further, it was demonstrated in Ye et al. (2024) that the purpose of temperature scaling or power transform is to enlarge the CMI values of temperature scaled (or power transformed) clusters, and enlarging CMI values in turn improves the performance of distilled students. Since here we want to achieve the opposite,

<sup>&</sup>lt;sup>1</sup>There are two reasons for us to adopt this definition of distillability. First, as shown theoretically in Zheng & YANG (2024), knowledge distillation reduces to label smoothing in the limit as the temperature approaches infinity. Therefore, if the teacher is distillable, the distilled student should perform no worse than the LS student. Second, if a knockoff student cannot outperform the LS student, then there is no incentive for the teacher to be leveraged since the LS student can be trained on its own.

we want to make sure that all CMI values of all power transformed clusters can be made small. To this end, we further extend the label specific CMI  $I(X; \hat{Y} | Y = y)$  and the CMI  $I(X; \hat{Y} | Y)$  to  $I(X; \hat{Y}^{\alpha} | Y = y)$  and  $I(X; \hat{Y}^{\alpha[Y]} | Y)$ , respectively, so that  $I(X; \hat{Y}^{\alpha} | Y = y)$  measures the concentration of the power transformed cluster corresponding to label y with power  $\alpha$ , and  $I(X; \hat{Y}^{\alpha[Y]} | Y)$  measures the average concentration across all power transformed clusters with power  $\alpha[Y]$ , where different clusters may be power transformed with different power  $\alpha$ . Notably, allowing separate temperature scaling (i.e., power transformation) for each class is essential, as an adaptive adversary may apply class-specific scaling to selectively alter cluster concentration and improve the knockoff student's performance.

Based on the above discussion and towards building undistillable DNNs, we then propose a new training method called CMI minimized method, which trains a DNN by jointly minimizing the CE loss and all CMI values of all power transformed clusters, i.e., jointly minimizing the CE loss and  $I(X; \hat{Y}^{\alpha[Y]} | Y), \forall \alpha[Y] > 0$ .

The resulting trained DNN is referred to as the CMI minimized (CMIM) DNN. The contributions of the paper are summarized as follows:

• An insight is provided that in order for a DNN to be undistillable, it is desirable for the DNN to possess the trait that each cluster of the DNN's output probability distributions corresponding to each label is highly concentrated to the extent that all probability distributions within the cluster more or less collapse into one probability distribution close to the one-hot probability vector of that label.

• We extend the label specific CMI  $I(X; \hat{Y} | Y = y)$  and the CMI  $I(X; \hat{Y} | Y)$  to  $I(X; \hat{Y}^{\alpha} | Y = y)$  and  $I(X; \hat{Y}^{\alpha[Y]} | Y)$ , respectively, so that  $I(X; \hat{Y}^{\alpha} | Y = y)$  measures the concentration of the power transformed cluster corresponding to label y with power  $\alpha$ , and  $I(X; \hat{Y}^{\alpha[Y]} | Y)$  measures the average concentration across all power transformed clusters with power  $\alpha[Y]$ , where different clusters may be power transformed with different power  $\alpha$ .

• We develop a novel training method dubbed CMI minimized method to train a DNN by jointly minimizing the CE loss and all CMI values of all power transformed clusters with the resulting trained DNN referred to as the CMIM DNN.

• To the best of our knowledge, our method is the first in the literature capable of training undistillable DNNs that remain robust against a wide range of KD methods. Furthermore, for the notion of undistillability, we are the first to employ the formal definition introduced in Yang & Ye (2024).

• We show, by extensive experiments over three popular image classification datasets, namely CIFAR-100 (Krizhevsky et al., 2012), TinyImageNet (Le & Yang, 2015) and ImageNet (Deng et al., 2009), that CMIM DNNs have very small CMI values and are indeed undistillable by all tested KD methods existing in the literature. That is, the knockoff students distilled by these KD methods from the CMIM models underperform the respective LS students. On the other hand, models trained by defense training methods proposed in the literature are all distillable.

• In addition, we show that the CMIM models achieve a higher classification accuracy compared to those trained with the conventional CE loss.

### 2 Related Works

In this section, we mention some defense methods against the threat posed by knockoff students attempting to steal the IP of pre-trained DNNs via logit-based KD methods. For a thorough review of related works, including detailed discussions about recent logit-based KD methods, please refer to Appendix B. These defense methods can be mainly categorized into two groups: (i) model stealing resistant training methods that specifically train DNNs to reduce the accuracy of knockoff students while maintaining the original classification accuracy of the model (Ma et al., 2021; Wang et al., 2022); and (ii) post-training defense methods that perform minimal perturbations to the pre-trained model's predictions to mislead the knockoff student (Lee et al., 2019; Orekondy et al., 2020; Cheng & Cheng, 2023). Nonetheless, in Section 5, we will show that models trained by all these defense methods are indeed distillable.

### 3 Notation and Preliminaries

#### 3.1 Notation

The set of real numbers is denoted by  $\mathbb{R}$ . Vectors are denoted by bold-face letters (e.g.,  $\boldsymbol{w}$ ). The *i*-th element of vector  $\boldsymbol{w}$  is denoted by  $\boldsymbol{w}[i]$ . For two vectors  $\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^C$ , the inequality  $\boldsymbol{u} \leq \boldsymbol{v}$  implies that  $\boldsymbol{u}[i] \leq \boldsymbol{v}[i]$ ,  $\forall i \in [C]$ . For a positive integer K, let  $[K] \triangleq \{1, \dots K\}$ . Assume that there are C class labels with [C] as the set of class labels. Let  $\mathcal{P}([C])$  denote the set of all C dimensional probability distributions. For any two probability distributions  $P_1, P_2 \in \mathcal{P}([C])$ , the CE and Kullback-Leibler (KL) divergence between  $P_1$  and  $P_2$ are denoted by  $H(P_1, P_2)$  and  $KL(P_1, P_2)$ , respectively. For any  $\boldsymbol{y} \in [C]$  and  $P \in \mathcal{P}([C])$ , write the CE of the one-hot probability distribution corresponding to  $\boldsymbol{y}$  and P as  $H(\boldsymbol{y}, P)$ .

For any differentiable function  $f(\cdot)$ ,  $\nabla_{\boldsymbol{w}} f(\cdot)$  denotes its gradient vector w.r.t. vector  $\boldsymbol{w}$ .

For any pair of random variables (X, Y), denote its joint probability distribution by  $P_{X,Y}(x, y)$  or simply P(x, y) whenever there is no ambiguity, the marginal distribution of Y by  $P_Y(y)$ , and the conditional distribution of Y given X = x by  $P_{Y|X}(\cdot|x)$ . The mutual information between two random variables X and Y is denoted by I(X, Y), and the CMI of X and Y given a third random variable Z is I(X, Y|Z).

We regard a classification DNN as a mapping from raw data  $x \in \mathbb{R}^d$  to a probability distribution  $q_x \in \mathcal{P}([C])$ . Given a DNN:  $x \in \mathbb{R}^d \to q_x$ , let  $\boldsymbol{\theta}$  denote its weight vector consisting of all its connection weights; whenever there is no ambiguity, we also write  $q_x$  as  $q_{x,\boldsymbol{\theta}}$ .

#### 3.2 Label Smoothing

Label smoothing (LS) (Pereyra et al., 2017) is a regularization technique that prevents peaked output probability distributions, leading to better generalization, by minimizing the objective function:

$$\mathcal{L}^{LS} = (1 - \epsilon)\mathsf{H}(y, q_x) + \epsilon\mathsf{H}(u, q_x),\tag{1}$$

where u is the uniform distribution over C classes, and  $\epsilon$  controls the strength of the regularization.

#### 3.3 Power Transform of Probability Distribution

In a "black-box" teacher setting, where only the output probability vectors (and not the logits) of the teacher are accessible to the public, applying temperature scaling directly over the logits of the teacher is not feasible in training knockoff students. In this case, KD training can resort to applying "power transformation of probability distribution" directly to the output probability vectors (Zheng & YANG, 2024). Specifically, given  $P \in \mathcal{P}([\mathcal{C}])$ , and a non-negative real number  $\alpha$ , the power transform of P is another probability distribution define as

$$P^{\alpha}[i] = \frac{(P[i])^{\alpha}}{\sum_{j \in [C]} (P[j])^{\alpha}}, \quad \forall i \in [C].$$

$$\tag{2}$$

It is not hard to verify that the power transformed probability distribution  $P^{\alpha}$  is equal to the softmax of the logits scaled by temperature  $T = 1/\alpha$ . Therefore, temperature scaling can be equivalently operated directly on the output probability distribution through power transform.

#### 3.4 CMI value of a DNN

As discussed in Yang et al. (2025), for a classifier  $f : x \in \mathbb{R}^d \to q_x$ , let  $\hat{Y}$  be the random label predicted by the f with probability  $q_X[\hat{Y}]$  in respond to the input X. For each cluster corresponding to label  $y \in [C]$ , we have

$$I(X; \hat{Y}|Y = y) = \sum_{x} P_{X|Y}(x|y) \left[ \sum_{i=1}^{C} P_{\hat{Y}|XY}(\hat{Y} = i|x, y) \ln \frac{P_{\hat{Y}|XY}(\hat{Y} = i|x, y)}{P_{\hat{Y}|Y}(\hat{Y} = i|Y = y)} \right]$$
(3)

$$= \mathbb{E}_{X|Y} \left[ \left( \sum_{i=1}^{C} q_X[i] \ln \frac{q_X[i]}{P_{\hat{Y}|Y}(\hat{Y} = i|Y = y)} \right) | Y = y \right] = \mathbb{E}_{X|Y} \left[ \text{KL} \left( q_X, s_y \right) | Y = y \right],$$
(4)

where  $P_{\hat{Y}|XY}(\hat{Y} = i|x, y) = q_x[i]$  follows from the Markov chain  $Y \to X \to \hat{Y}$ , and  $s_y = P_{\hat{Y}|Y}(\cdot|y) = \mathbb{E}_{X|Y}[q_X|Y = y]$ . I $(X; \hat{Y}|Y = y)$  measures the concentration of the cluster corresponding to label  $y \in [C]$ . Averaging over all clusters corresponding to all labels y, we get

$$I(X; \hat{Y}|Y) = \sum_{y \in [C]} P_Y(y) I(X; \hat{Y}|Y=y) = \mathbb{E}_{XY} \left[ KL\left(q_X, s_Y\right) \right].$$
(5)

 $I(X; \hat{Y}|Y)$  measures the average concentration across all clusters.

When the distribution  $P_{X,Y}$  is unknown, we can approximate the CMI of f by its empirical value from a data sample (a training dataset or mini-batch thereof)  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^m$ . To this end, let  $\mathcal{D}_y = \{1 \leq j \leq m : y_j = y\}$ . Denote the size of  $\mathcal{D}_y$  by  $|\mathcal{D}_y|$ . The empirical values of each label specific CMI and the CMI can be calculated as follows

$$\mathbf{I}^{emp}(X;\hat{Y}|Y=y) = \frac{1}{|\mathcal{D}_y|} \sum_{i \in \mathcal{D}_y} \mathrm{KL}(q_{x_i}, s_y^{emp}),\tag{6}$$

$$I^{emp}(X; \hat{Y}|Y) = \frac{1}{m} \sum_{i=1}^{m} KL(q_{x_i}, s_{y_i}^{emp}),$$
(7)

where 
$$s_y^{emp} = \frac{1}{|\mathcal{D}_y|} \sum_{i \in \mathcal{D}_y} q_{x_i}, \forall y \in [C].$$
 (8)

### 4 CMI Minimized Method

In this section, we present our CMI minimized method. We begin with extending  $I(X; \hat{Y}|Y = y)$  and  $I(X; \hat{Y}|Y)$  to the case of power transformed clusters.

#### 4.1 Information Quantities for Power Transformed Clusters

Consider a classification DNN:  $f: x \in \mathbb{R}^d \to q_x$  which maps input sample instances x with different labels into clusters of probability distributions  $q_x$  in the space  $\mathcal{P}([C])$ , with one cluster per label. For each label  $y \in [C]$ , apply the power transform with power  $\alpha$  to each probability distribution  $q_x$  within the cluster corresponding to the label y. Then, we obtain a power transformed cluster. To measure the concentration of the power transformed cluster, we extend  $I(X; \hat{Y}|Y = y)$  to the following information quantity

$$I(X; \hat{Y}^{\alpha} | Y = y) = \mathbb{E}_{X|Y} \left[ KL\left(q_X^{\alpha}, s_{y,\alpha}\right) | Y = y \right], \tag{9}$$

where  $s_{y,\alpha} = \mathbb{E}_{X|Y}[q_X^{\alpha}|Y=y]$ . Note that if we regard  $\hat{Y}^{\alpha}$  as the random label predicted by f with probability  $q_X^{\alpha}(\hat{Y}^{\alpha})$  in response to the input sample X, i.e., given X,  $\hat{Y}^{\alpha}$  is equal to a label c with probability  $q_X^{\alpha}(c)$ ,  $\forall c \in [C]$ , then  $I(X; \hat{Y}^{\alpha}|Y=y)$  is exactly the CMI between X and  $\hat{Y}^{\alpha}$  given Y = y. Thus,  $I(X; \hat{Y}^{\alpha}|Y=y)$  measures the concentration of the power transformed cluster corresponding to y.

Now, we go one step further and allow different clusters to be power transformed with different powers. Suppose that the cluster corresponding to label y is power transformed with power  $\alpha[y]$ . Let  $\hat{Y}^{\alpha[Y]}$  be the random label predicted by f with probability  $q_X^{\alpha[Y]}(\hat{Y}^{\alpha[Y]})$  in response to the input sample X given Y. That is, given Y = y and X = x,  $\hat{Y}^{\alpha[Y]}$  is equal to c with probability  $q_x^{\alpha[y]}(c)$  for any  $c \in [C]$ . We can then extend  $I(X; \hat{Y}|Y)$  to  $I(X; \hat{Y}^{\alpha[Y]}|Y)$ 

$$I(X; \hat{Y}^{\alpha[Y]}|Y) = \mathbb{E}_{XY} \left[ KL\left(q_X^{\alpha[Y]}, s_{Y,\alpha[Y]}\right) \right],$$
(10)

$$= \sum_{y \in [C]} P_Y(y) \left[ \mathbb{E}_{X|Y} \left[ \mathrm{KL} \left( q_X^{\alpha[y]}, s_{y,\alpha[y]} \right) | Y = y \right] \right]$$
(11)

$$= \sum_{y \in [C]} P_Y(y) I(X; \hat{Y}^{\alpha[y]} | Y = y),$$
(12)

where for each  $y \in [C]$ ,

$$s_{y,\alpha[y]} = P_{\hat{Y}^{\alpha[Y]}|Y}(\cdot|y) = \sum_{x} P_{X|Y}(x|y)q_x^{\alpha[y]} = \mathbb{E}_{X|Y}\left[q_X^{\alpha[y]}|Y=y\right].$$
(13)

Note that  $I(X; \hat{Y}^{\alpha[Y]}|Y)$  is exactly the CMI between X and  $\hat{Y}^{\alpha[Y]}$  given Y and measures the average concentration across all power transformed clusters with power function  $\alpha[Y]$ . However, Y, X, and  $\hat{Y}^{\alpha[Y]}$  do not form a Markov chain anymore.

When the distribution  $P_{X,Y}$  is unknown, we can approximate  $I(X; \hat{Y}^{\alpha[Y]}|Y = y)$  and  $I(X; \hat{Y}^{\alpha[Y]}|Y)$  by their respective empirical values from a data sample (a training dataset or mini-batch thereof)  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{m}$ :

$$\mathbf{I}^{emp}(X; \hat{Y}^{\alpha[Y]} | Y = y) = \frac{1}{|\mathcal{D}_y|} \sum_{i \in \mathcal{D}_y} \mathrm{KL}(q_{x_i}^{\alpha[y]}, s_{y, \alpha[y]}^{emp}),$$
(14)

$$I^{emp}(X; \hat{Y}^{\alpha[Y]}|Y) = \frac{1}{m} \sum_{i=1}^{m} \text{KL}(q_{x_i}^{\alpha[y_i]}, s_{y_i, \alpha[y_i]}^{emp}),$$
(15)

where 
$$s_{y,\alpha[y]}^{emp} = \frac{1}{|\mathcal{D}_y|} \sum_{i \in \mathcal{D}_y} q_{x_i}^{\alpha[y]}, \forall y \in [C].$$
 (16)

As discussed in Section 1, an undistillable DNN should exhibit the trait that each of these clusters is highly concentrated and ideally collapses into a single probability distribution that closely resembles the one-hot probability vector for that label.

**Remark 1.** While we leverage the concept of CMI from Yang et al. (2025), the way it is calculated in our work significantly differs from how it is calculated in Yang et al. (2025). In Yang et al. (2025),  $I(X; \hat{Y}|Y)$  is calculated under the assumption of a Markov chain  $Y \to X \to \hat{Y}$ . In contrast, we quantify cluster compactness using  $I(X; \hat{Y}^{\alpha[Y]}|Y)$ , where  $\hat{Y}^{\alpha[Y]}$  explicitly depends on Y, violating the Markov assumption. Moreover, minimizing  $I(X; \hat{Y}^{\alpha[Y]}|Y)$  over all possible values of  $\alpha$  introduces additional challenges, which we address in the next subsection.

#### 4.2 Framework for Minimizing CMI Values of Power Transformed Clusters

Towards building an undistillable DNN, we now train a DNN  $f : x \in \mathbb{R}^d \to q_x$  by jointly minimizing the CE loss and all CMI values of all power transformed clusters. Let

$$\boldsymbol{\alpha} = [\alpha[1], \alpha[2], \dots, \alpha[C]],$$

and write each  $q_x$  as  $q_{x,\theta}$ . In our CMI minimized method, the objective function we want to minimize is

 $\mathbb{E}_{XY} | \mathsf{H}$ 

$$(Y, q_{X, \theta}) \Big] + \lambda \max_{\alpha} I(X; \hat{Y}^{\alpha[Y]} | Y), \tag{17}$$



Figure CMI The value 1:  $I(X; \hat{Y}^{\alpha[Y]}|Y$ y)for three = classes randomly selected  $y = \{\text{beaver, orchids, motorcycle}\}$ vs the power transform factor  $\alpha$ ; the model is ResNet-50 pre-trained on CIFAR-100. As observed, as  $\alpha$  grows, the CMI value becomes larger, peaks, and then gradually becomes smaller.

where  $\lambda > 0$  is a hyper-parameter trading the CE loss with the maximum CMI, and the maximization over  $\boldsymbol{\alpha}$  is taken over the region  $0 \le \alpha[i] \le \beta$ ,  $1 \le i \le C$ . The optimization problem then becomes

$$\min_{\boldsymbol{\theta}} \left\{ \mathbb{E}_{XY} \Big[ \mathsf{H}(Y, q_{X, \boldsymbol{\theta}}) \Big] + \lambda \max_{\boldsymbol{\alpha}} \mathrm{I}(X; \hat{Y}^{\boldsymbol{\alpha}[Y]} | Y) \right\} \\
= \min_{\boldsymbol{\theta}} \left\{ \mathbb{E}_{XY} \Big[ \mathsf{H}(Y, q_{X, \boldsymbol{\theta}}) \Big] + \lambda \max_{\boldsymbol{\alpha}} \sum_{y} P_{Y}[y] \mathrm{I}(X; \hat{Y}^{\boldsymbol{\alpha}[y]} | Y = y) \right\} \tag{18}$$

$$= \min_{\boldsymbol{\theta}} \Big\{ \mathbb{E}_{XY} \Big[ \mathsf{H}(Y, q_{X, \boldsymbol{\theta}}) \Big] + \lambda \sum_{y} P_{Y}[y] \max_{\alpha[y]} \mathrm{I}(X; \hat{Y}^{\alpha[y]} | Y = y) \Big\}.$$
(19)

In order to get a better understanding about the behavior of the second term in the objective function of equation 18 w.r.t.  $\alpha$ , we depict in Figure 1 I(X;  $\hat{Y}^{\alpha[Y]}|Y = y)$  vs  $\alpha[y]$  for three randomly-selected classes y using a pre-trained ResNet-50 on CIFAR-100. In Figure 1,  $\max_{\alpha} I(X; \hat{Y}^{\alpha}|Y = y)$  is achieved at a value of  $\alpha$  which is between 0.25 and 0.75. In Theorem 2 of Appendix D, we further show that for each label y,  $I(X; \hat{Y}^{\alpha}|Y = y)$  as a function of  $\alpha$  is continuously differentiable.

However, finding an algorithmic solution to the min-max problem in equation 18 to equation 19 is challenging. To overcome this difficulty, we next develop a more tractable expression for  $\max_{\alpha} I(X; \hat{Y}^{\alpha}|Y = y)$ . At this point, we invoke the following theorem, which will be proved in Appendix E.

**Theorem 1.** For any label y,

$$\max_{\alpha} I(X; \hat{Y}^{\alpha} | Y = y) = \lim_{\omega \to \infty} \frac{1}{\omega} \ln \frac{1}{\beta} \int_{0}^{\beta} \exp\left\{\omega I(X; \hat{Y}^{\alpha} | Y = y)\right\} d\alpha.$$
(20)

Therefore, when  $\omega$  is large,  $\max_{\alpha} I(X; \hat{Y}^{\alpha}|Y = y)$  can be approximated by

$$\max_{\alpha} I(X; \hat{Y}^{\alpha} | Y = y) \approx \frac{1}{\omega} \ln \frac{1}{\beta} \int_{0}^{\beta} \exp\left\{\omega I(X; \hat{Y}^{\alpha} | Y = y)\right\} d\alpha$$
(21)

$$\approx \frac{1}{\omega} \ln \left[ \frac{1}{N} \sum_{i=1}^{N} \exp \left\{ \omega \mathbf{I}(X; \hat{Y}^{\alpha_i} | Y = y) \right\} \right], \tag{22}$$

where N is relatively large, and  $\alpha_i = i\beta/N$ .

Now plugging equation 22 into equation 19, we have

$$\min_{\boldsymbol{\theta}} \left\{ \mathbb{E}_{XY} \Big[ \mathsf{H}(Y, q_{X, \boldsymbol{\theta}}) \Big] + \frac{\lambda}{\omega} \sum_{y} P_{Y}[y] \ln \left[ \frac{1}{N} \sum_{i=1}^{N} \exp \left\{ \omega \mathrm{I}(X; \hat{Y}^{\alpha_{i}} | Y = y) \right\} \right] \right\}.$$
(23)

Note that the second term in the objective function of equation 23 is not amenable to parallel computation via GPU due to the dependency of KL divergence on  $s_{y,\alpha_i}$ , the centroid of the power transformed cluster corresponding to Y = y with power  $\alpha_i$ . To get around this difficulty, we follow the approach in Yang et al. (2025) and introduce dummy distributions  $Q_{y,i} \in \mathcal{P}([C])$  for each (y,i) to rewrite  $I(X; \hat{Y}^{\alpha_i}|Y = y)$  as follows

$$I(X; \hat{Y}^{\alpha_i} | Y = y) = \mathbb{E}_{X|Y} \left[ \mathrm{KL} \left( q_{X, \boldsymbol{\theta}}^{\alpha_i}, s_{y, \alpha_i} \right) \mid Y = y \right]$$
$$= \min_{Q_{y, i}} \mathbb{E}_{X|Y} \left[ \mathrm{KL} \left( q_{X, \boldsymbol{\theta}}^{\alpha_i}, Q_{y, i} \right) \mid Y = y \right],$$
(24)

where the minimum in the above is achieved when

$$Q_{y,i} = s_{y,\alpha_i} = \mathbb{E}_{X|Y} \left[ q_{X,\boldsymbol{\theta}}^{\alpha_i} | Y = y \right].$$
(25)

Combining equation 24 with equation 23, we are led to solve the double minimization problem  $^2$ 

$$\min_{\boldsymbol{\theta}} \left\{ \mathbb{E}_{XY} \Big[ \mathsf{H}(Y, q_{X, \boldsymbol{\theta}}) \Big] + \frac{\lambda}{\omega} \sum_{y} P_{Y}[y] \ln \left[ \frac{1}{N} \sum_{i=1}^{N} \exp \left\{ \omega \min_{Q_{y, i}} \mathbb{E}_{X|Y} \Big[ \mathrm{KL} \left( q_{X, \boldsymbol{\theta}}^{\alpha_{i}}, Q_{y, i} \right) \mid Y = y \Big] \right\} \right] \right\}$$
(26)

$$= \min_{\boldsymbol{\theta}} \min_{\{Q_{y,i}\}_{y \in [C], i \in [N]}} \left\{ \mathbb{E}_{XY} \left[ \mathsf{H}(Y, q_{X, \boldsymbol{\theta}}) \right] + \frac{\lambda}{\omega} \sum_{y} P_{Y}[y] \ln \left[ \frac{1}{N} \sum_{i=1}^{N} \exp \left\{ \omega \mathbb{E}_{X|Y} \left[ \mathrm{KL} \left( q_{X, \boldsymbol{\theta}}^{\alpha_{i}}, Q_{y, i} \right) \mid Y = y \right] \right\} \right] \right\}$$
(27)

When the distribution  $P_{X,Y}$  is unknown, it can be approximated by its empirical distribution from a data sample (a training dataset or mini-batch thereof)  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^m$ . The objective function in the double minimization equation 27 then becomes

$$J_{\mathcal{D}}(\boldsymbol{\theta}, \{Q_{y,i}\}_{y \in [C], i \in [N]}) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \mathsf{H}(y, q_{x,\boldsymbol{\theta}}) + \frac{\lambda}{\omega} \sum_{y} \frac{|\mathcal{D}_{y}|}{|\mathcal{D}|} \ln \left[ \frac{1}{N} \sum_{i=1}^{N} \exp\left\{ \frac{\omega}{|\mathcal{D}_{y}|} \sum_{j \in \mathcal{D}_{y}} \mathrm{KL}\left(q_{X_{j},\boldsymbol{\theta}}^{\alpha_{i}}, Q_{y,i}\right) \right\} \right].$$
(28)

#### 4.3 Algorithm for Solving the Optimization in equation 27

The double minimization optimization problem in equation 27 naturally lends us an alternating algorithm that optimizes  $\theta$  and  $\{Q_{y,i}\}_{y \in [C], i \in [N]}$  alternatively to minimize the objective function in equation 27 or Equation (28), given the other is fixed.

Given  $\{Q_{y,i}\}_{y\in[C],i\in[N]}$ ,  $\boldsymbol{\theta}$  can be updated using the same first-order optimization method as in conventional deep learning, such as stochastic gradient descent applied over mini-batches.

Following Yang et al. (2025), given  $\theta$ , for each class y,  $\{Q_{y,i}\}_{i \in [N]}$  can be updated according to equation 25 in the following manner: (1) we randomly sample a mini-batch of samples  $|\mathfrak{B}_y|$  instances from the training set with ground truth label y; (2)  $\{Q_{y,i}\}_{i \in [N]}$  can be updated as

$$Q_{y,i} = \frac{\sum_{x \in \mathfrak{B}_y} q_{x,\boldsymbol{\theta}}^{\alpha_i}}{|\mathfrak{B}_y|} \quad \forall i \in [N].$$
<sup>(29)</sup>

The proposed alternating algorithm for optimization problem equation 27 is summarized in Algorithm 1<sup>3</sup>. To simplify our notation, we use  $(\cdot)_b^t$  to indicate parameters at the *b*-th batch updation during the *t*-th alternating iteration of the algorithm. We further write  $(\cdot)_B^t$  as  $(\cdot)^t$  whenever needed, set  $(\cdot)_0^t = (\cdot)^{t-1}$ .

**Remark 2.** The compactness of output clusters alone is insufficient to ensure an undistillable DNN. Undistillability is a significantly stronger property. For instance, LS can improve the compactness of the feature space of a DNN, which may make the output probability of the DNN more compact (Müller et al., 2019). However, this compactness is not enough to make the DNN undistillable (see Section 5)

### 5 Experiments

In this section, we demonstrate the effectiveness of CMIM by comparing it with several state-of-the-art alternatives. Specifically, we first report the accuracy that a knockoff student can achieve by deploying

 $<sup>^{2}</sup>$ In practice, solving the double minimization problem introduces only minor runtime overhead, as the inner optimization problem has an analytic solution and can be parallelized efficiently on modern GPUs. To demonstrate the efficiency of CMIC, we report the wall-clock training time and compare the computational overhead of CMIM and CE in Appendix J.

<sup>&</sup>lt;sup>3</sup>If the impact of the random mini-batch sampling and stochastic gradient descent is ignored, the alternating algorithm is guaranteed to converge in theory since given  $\theta$ , the optimal  $\{Q_{y,i}\}_{y \in [C], i \in [N]}$  can be found analytically via equation 29, although it may not converge to a global minimum.

Algorithm 1: Conditional Mutual Information Minimized (CMIM) Method. Input: Training set  $\mathcal{T}$ , mini-batches  $\{\mathcal{B}_b\}_{b\in[B]}$ , number of epochs T,  $\lambda$ ,  $\beta$ ,  $\omega$ , NInitialization: Initialize  $\theta^0$  and  $Q_{y,i}^0|_{y\in[C],i\in[N]}$ . for t = 1 to T do [Sampling  $\alpha_i$ ] Randomly select N samples  $\{\alpha_i\}_{i\in[N]}$  from interval  $[0,\beta]$ . for b = 1 to B do [Updating  $Q_{y,i}$ ] For each class y, construct mini-batch  $\{\mathcal{B}_y\}_{y\in[C]}$ . Update  $Q_{y,i}^t, \forall y \in [C]; \forall i \in [N],$ according to Equation (29). [Updating  $\theta$ ] Fix  $Q_{y,i}^t|_{y\in[C],i\in[N]}$ . Update  $\theta_{b-1}^t$  to  $\theta_b^t$  by stochastic gradient descent over the objective function 28. end end Output: Global model  $\theta^T$ .

different logit-based KD (attack) methods in Section 5.1. In all the experiments, when testing the distillibility of the trained DNNs using the benchmark defense methods and CMIM, we compare the knockoff student's accuracy (i) when it attempts to steal the IP of protected DNN using logit-based (attack) methods with (ii) when it trains its model using the LS. If the former outperforms the latter, we conclude that the knockoff makes the underlying DNN distillable. Next, in Section 5.2, we report the classification accuracy of the *protected* models trained by the different defense methods. Lastly, in Section 5.3, we visualize the output cluster of models trained by CMIM, CE and NT.

### 5.1 Knockoff Student Accuracy

• **Datasets:** We conduct extensive experiments on three image classification dataset, namely CIFAR-100 (Krizhevsky et al., 2012) TinyImageNet (Le & Yang, 2015) and ImageNet (Deng et al., 2009). For description of each dataset, please refer to Appendix F.

• Models: To show the effectiveness of CMIM, we use different model architectural families for teacher and knockoff student models. To this end, we pick models from VGG family (Simonyan & Zisserman, 2015), ResNet family (He et al., 2016) (shortened as RN), ShuffleNetV2 (Ma et al., 2018), shortened as SNV2, and Mobilenetv2 (Sandler et al., 2018) shortened as MNV2. Particularly, we have conducted experiments on the following (teacher-student) pairs for each dataset: (i) for CIFAR-100, we use four pairs {(VGG16-VGG11), (VGG16-SNV2), (RN50-VGG11), (RN50-RN18)}; (ii) for TinyImageNet, we use two pairs {(RN34-RN18), (RN50-SNV2)}; and for ImageNet we use two pairs {(RN34-RN18),(RN34-MNV2)}.

• Defense benchmark methods: For comprehensive comparisons, we benchmark CMIM with seven recently published defense methods: MAD (Orekondy et al., 2020), APGP (Cheng & Cheng, 2023), RSP (Lee et al., 2019), ST (Ma et al., 2022), NT (Ma et al., 2021), SNT (Wang et al., 2022), and LS<sup>4</sup> (Müller et al., 2019).

• Logit-based KD (attack) methods: We use three logit-based KD methods that are primarily designed for when the teacher-student models are in cooperating mode, namely KD (Hinton et al., 2015), DKD (Zhao et al., 2022), DIST (Huang et al., 2022a); and four logit-based KD attacks methods that a knockoff student can deploy to make the protected DNNs possibly distillable, namely MKD (Yang & Ye, 2024), HTC (Jandial et al., 2022), AVG (Keser & Toreyin, 2023), Knockoff (Orekondy et al., 2019). We report all the training setups, including all the hyper-parameters used for both defense and attack methods in Appendix G.1.

• **Results:** The accuracy that a knockoff student can achieve using the various (*defense-attack*) combinations is summarized in Table 1. For accuracy variances, please refer to Appendix I. In the table, we use the notation "K-student" to denote a knockoff student. The numbers in the column titled "Best" represent the highest

 $<sup>^{4}</sup>$ Although LS is not a defense method per se, it is observed that the models trained by LS reduce the knockoff student's accuracy. We discuss the rationale behind this in Appendix C.

Table 1: Top-1 accuracy (%) of the knockoff student on CIFAR-100, TinyImageNet and ImageNet dataset (the results for CIFAR-100 and TinyImageNet are averaged over 3 runs). Green upward arrows ( $\uparrow$ ) and red downward arrows ( $\downarrow$ ) indicate whether the knockoff student was able to render the underlying DNN distillable.

CIFAR-100											
Defense	Model	K-student	LS	KD	MKD	DKD	DIST	HTC	AVG	Knockoff	Best
	VOOIC	VGG11	71.94	68.55	72.08	53.32	69.21	71.19	70.03 🗸	61.44	72.08
MAD	VGG16	SNV2	72.65	$72.50 \downarrow$	$72.46 \downarrow$	7.64 \downarrow	$69.91 \downarrow$	71.37	72.86	70.87 \downarrow	72.86
MAD	DNFO	VGG11	71.94	72.00	72.04	54.29	71.57	70.76	70.73	61.73	72.04
	KN50	RN18	78.76	77.76 🗼	78.79	43.73 🗸	73.76 \downarrow	77.89 🗸	78.61	73.92	78.79
	VOOIC	VGG11	71.94	71.92	72.27	27.24	69.25	70.08	72.01	45.98	72.27 ↑
ADCD	VGG10	SNV2	72.65	73.10	73.75	$12.52 \downarrow$	71.04	$71.66 \downarrow$	73.20	9.48 👃	73.75 ↑
AFGF	DNEO	VGG11	71.94	71.91 🗸	72.11	9.74 🗸	$69.48 \downarrow$	$71.36 \downarrow$	71.92	34.71 🗸	72.11
	10100	RN18	78.76	$78.04 \downarrow$	$79.06 \uparrow$	$62.71 \downarrow$	$77.32 \downarrow$	77.82 🗸	77.90 \downarrow	2.57 👃	79.06 ↑
	VCC16	VGG11	71.94	71.42 🗸	$72.04 \uparrow$	70.22 🗸	70.80 \downarrow	70.40 🗸	$71.56 \downarrow$	31.04 🗸	72.04 ↑
RSP	19910	SNV2	72.65	$73.55 \uparrow$	$72.95 \uparrow$	$67.45 \downarrow$	$72.19 \downarrow$	$71.46 \downarrow$	72.27 \downarrow	26.09 🗸	73.55 ↑
1651	RN50	VGG11	71.94	$71.97 \uparrow$	$72.01 \uparrow$	$69.53 \downarrow$	$72.18 \uparrow$	70.87 \downarrow	$70.85 \downarrow$	$46.68 \downarrow$	72.18 ↑
	111100	RN18	78.76	77.78 🗸	77.79 \downarrow	77.01 🗸	$78.88 \uparrow$	78.00 🗸	78.13 🔶	$55.86 \downarrow$	78.88 ↑
VCC16	VGG11	71.94	$71.40 \downarrow$	$73.44 \uparrow$	71.47 🗸	$71.33 \downarrow$	70.77 \downarrow	$71.58 \downarrow$	$63.56 \downarrow$	73.44 ↑	
NT	10010	SNV2	72.65	72.44 🗸	$72.70 \uparrow$	6.24	72.04 🗸	70.75 🗸	$72.83 \uparrow$	6.32 👃	72.83 ↑
	BN50	VGG11	71.94	$72.01 \uparrow$	$72.03 \uparrow$	$71.55 \downarrow$	71.88 🗸	$70.16 \downarrow$	$71.94 \downarrow$	$62.94 \downarrow$	72.03 ↑
	10100	RN18	78.76	78.41 🗸	$78.92 \uparrow$	$79.26 \uparrow$	$78.99 \uparrow$	77.94 🗸	78.33 🗸	68.96 \downarrow	79.26 ↑
	VGG16	VGG11	71.94	$72.06 \uparrow$	$72.28 \uparrow$	4.92 👃	$71.98 \uparrow$	70.60 🗸	$71.63 \downarrow$	64.08 🗸	72.06 ↑
SNT	, 44410	SNV2	72.65	$72.94 \uparrow$	73.17	$72.78 \uparrow$	72.22 \downarrow	71.22 🗸	$72.74 \uparrow$	6.22	73.17 ↑
BN50	VGG11	71.94	$72.02 \uparrow$	72.12	72.32	71.70 🗸	70.66 \downarrow	$71.65 \downarrow$	$62.94 \downarrow$	72.32 ↑	
		RN18	78.76	78.25 🗸	78.48	78.82 ↑	78.14	78.45 🗸	78.38 🗸	67.71 \downarrow	78.82 ↑
VGG16	VGG11	71.94	$72.09 \uparrow$	$72.01 \uparrow$	$71.63 \downarrow$	$71.93 \downarrow$	$71.16 \downarrow$	71.63	$63.32 \downarrow$	72.09	
ST	ST	SNV2	72.65	72.64	72.67 ↑	70.53	72.24	71.32	72.42	69.46	72.67
RN50	VGG11	71.94	$72.00 \uparrow$	72.13	71.62	71.76	70.54	71.73	65.43	72.13	
		RN18	78.76	78.96 ↑	79.02	78.35	78.31	78.36	78.81	72.87↓	79.02
	VGG16	VGGII	71.94	71.90	72.00 ↑	71.57	70.89	70.66	71.76	$63.49 \downarrow$	72.00
LS		SNV2	72.65	72.87	$73.52\uparrow$	70.01	71.49	71.70	73.01	65.20	73.52
	RN50	VGGII	71.94	77.70	71.99	71.95	70.77	70.80	71.88	$62.29 \downarrow$	71.99
		KN18 VCC11	71.04	71.07	71.64	71.56	$\frac{18.33}{70.24}$	71 71	71.49	$03.30 \downarrow$	79.37
	VGG16	VGGII	71.94	(1.8(↓ 70.52↓	$(1.04 \downarrow 71.44 \downarrow$	$(1.30 \downarrow 72.46 \downarrow$	70.34	(1.(1 ↓ 71 E0 ↓	(1.42 ↓ 71.04 ↓	$60.89 \downarrow$	$(1.8) \downarrow$
CMIM		SINV2 VCC11	72.03	$71.53 \downarrow$	$71.44 \downarrow$	71.77	71.40	60.22	71.94	$04.40 \downarrow$	71.00
	RN50	RN18	71.94	$71.04 \downarrow$ 78.01	71.04 J	78 13	$71.00 \downarrow$ 77.56	$09.34 \downarrow$	78.64	65.88	78.64
		minio	18.10	10.21	78.10 ↓	10.13	11.50 ↓	11.20 🗸	10.04 🗸	05.88 4	10.04 🗸
					TinyIm	ageNet					
RSP	RN34	RN18	63.56	$63.54 \downarrow$	$64.32$ $\uparrow$	$64.01 \uparrow$	$63.27 \downarrow$	$63.54 \downarrow$	$62.15 \downarrow$	$55.43 \downarrow$	64.32 ↑
1001	RN50	SNV2	60.61	$60.18 \downarrow$	$60.76 \uparrow$	$56.26 \downarrow$	$56.43 \downarrow$	$60.96 \uparrow$	$60.15 \downarrow$	$54.01 \downarrow$	60.96
ST	RN34	RN18	63.56	$63.96 \uparrow$	$64.12$ $\uparrow$	$63.25 \downarrow$	$63.51 \downarrow$	$63.49 \downarrow$	$63.84 \uparrow$	$57.42 \downarrow$	64.12 ↑
	RN50	SNV2	60.61	$61.23 \uparrow$	$61.36 \uparrow$	$60.43 \downarrow$	$60.32 \downarrow$	60.22 🗸	$61.13 \uparrow$	55.84 \downarrow	61.36 ↑
NT	RN34	RN18	63.56	$63.27 \downarrow$	$64.49 \uparrow$	$64.67 \uparrow$	$63.43 \downarrow$	$63.50 \downarrow$	$64.43 \uparrow$	$53.11 \downarrow$	64.67 ↑
	RN50	SNV2	60.61	$59.57 \downarrow$	$61.55 \uparrow$	$31.55 \downarrow$	$60.03 \downarrow$	$60.98 \uparrow$	$60.31 \downarrow$	$50.94 \downarrow$	$61.55 \uparrow$
LS	RN34	RN18	63.56	$63.74 \uparrow$	$64.01 \uparrow$	$64.23 \uparrow$	$63.51 \downarrow$	$64.20 \uparrow$	63.04	57.43 \downarrow	64.23 ↑
	RN50	SNV2	60.61	$60.32 \downarrow$	$60.93 \uparrow$	$60.74 \uparrow$	60.11 \downarrow	$60.46 \downarrow$	$60.14 \downarrow$	$52.96 \downarrow$	60.93 ↑
CMIM	RN34	RN18	63.53	62.89 🗸	63.15 🗸	62.94	63.28 🗸	61.57	62.96 🗸	56.13	63.28
	RN50	SNV2	60.61	57.57 🗸	59.32 \downarrow	60.58	59.41 \downarrow	59.33 🗸	$60.42 \downarrow$	$56.91 \downarrow$	60.58 🗸
					Imag	geNet					
am	DATE	RN18	70.89	70.74 👃	71.02 ↑	70.02	69.94 J	70.91 ↑	71.00 ↑	63.24	71.02 ↑
ST	RN34	MNV2	70.93	71.03	71.25	69.32	70.53	70.69	71.06	54.53 J	71.25
(D) (D) (	DNG	RN18	70.89	70.44	70.69	69.97	70.59	70.63	70.53 👃	59.34	70.69
CMIM	KN34	MNV2	70.93	70.21	70.72	69.97	70.44	70.86	70.20	55.24	70.86

accuracy obtained for each respective row, indicating the best possible performance a knockoff student can achieve using any of the listed distillation methods.

As observed in Table 1, regardless of whether the teacher-student architectures are the same or different, DNNs trained with CMIM remain undistillable across all distillation methods. This is in stark contrast to DNNs trained using other defense techniques, which can still be successfully distilled to a certain degree. The results indicate that prior defense strategies do not offer complete resistance against knockoff students, whereas CMIM effectively prevents distillation, making it significantly more robust in protecting model knowledge.

We also perform an ablation study on the hyperparameters of CMIM—namely  $\beta$ , N, and  $\omega$ —as detailed in Appendix K. The key findings are as follows:

- Effect of  $\beta$ : The accuracy of the knockoff student drops sharply when  $\beta \geq 1$ , suggesting that large values may destabilize optimization or excessively penalize the CMI term. This highlights the importance of carefully tuning  $\beta$  to effectively balance the cross-entropy and CMI objectives.
- Effect of the number of samples N: Increasing the number of power samples N leads to a monotonic decrease in the knockoff student's accuracy. This indicates that moderate values of N are sufficient to capture the necessary diversity for robust CMI estimation.
- Effect of the power coefficient  $\omega$ : Setting  $\omega > 30$  can lead to numerical instability, resulting in NaNs due to excessive exponentiation. Interestingly, the worst knockoff student performance is observed at  $\omega = 20$  and  $\omega = 30$ , suggesting that these settings best approximate the extreme concentration behavior of  $\omega \to \infty$ , thereby enhancing undistillability.

Additionally, we investigate CMIC's impact on model calibration in Appendix M.

### 5.2 Accuracy of Protected Models

In this section, we report the top-1 accuracy of the *protected* models in Table 1 trained using the benchmark defense methods with those trained by CMIM. The results are summarized in Tables 2 and 3. As observed, the models trained by CMIM have the highest classification accuracy compared to the benchmark methods. This is because for the models trained by CMIM, the clusters corresponding to the output probability of the DNNs are very concentrated, facilitating easier classification of samples from different classes.

Table 2: Top-1 accuracy (%) of models trained by defense methods on CIFAR-100 and TinyImageNet. The best and second best results are **bolded** and <u>underlined</u>, respectively.

CIFAR100								TinyImageNet								
Model	CE	MAD	APGP	RSP	ST	NT	SNT	LS	CMIM	Model	CE	RSP	ST	NT	LS	CMIM
VGG16	73.75	73.75	73.84	73.71	73.75	73.75	72.59	73.90	73.84	RN34	65.39	65.21	65.39	65.23	65.45	65.99
RN50	77.81	77.81	77.56	77.63	77.81	77.31	77.77	78.45	78.72	RN50	<u>66.14</u>	65.91	66.13	66.06	66.09	66.93

Table 3: Top-1 accuracy (%) of models trained by defense methods on ImageNet.

ImageNet									
Model	CE	ST	CMIM						
RN34	<u>73.31</u>	73.30	73.69						

The results in Table 2 motivate us to test the top-1 accuracy of additional models trained by CMIM and compare them with those trained by CE loss (see Appendix H).

It is worth noting that the primary focus of this paper is not on increasing the accuracy of DNNs but on developing a method to train undistillable DNNs. While many existing methods in the literature can enhance a DNN's accuracy, they do not address the critical challenge of making DNNs undistillable.



Figure 2: Visualization of three projected probability clusters for ResNet-50 trained on CIFAR-100 using (a) LS, (b) NT, and (c) CMIM.

Our approach is the first in the literature that effectively trains undistillable models robust against a wide range of existing KD methods. The improvement in accuracy observed in our results is a by-product of our method and not its primary goal. This improvement arises from the unique properties of our approach rather than replicating the effects of label smoothing or similar techniques.

### 5.3 Visualizing the Output Clusters

In this subsection, we aim to visualize the output clusters for models trained using CE, NT, and CMIM. To achieve this, we follow the visualization approach introduced by Yang et al. (2025). Specifically, we randomly select three labels from the CIFAR-100 dataset. For each probability distribution corresponding to these three labels, we extract only the probabilities associated with these selected labels and normalize them to form three-dimensional probability vectors. These vectors are then projected onto a two-dimensional simplex, allowing us to visualize the clustering behavior of each model. By applying this transformation, we obtain a clear representation of how the models distribute their probability mass across different output categories.

The resulting simplexes for ResNet-50 models trained with LS, Nasty Teacher, and the CMIM framework are shown in Figure 6<sup>5</sup>. To ensure a consistent comparison, we applied the same power transform  $\alpha = 4$  for all visualizations. As observed, the clusters for the model trained with CMIM are highly concentrated near the corners of the simplex, closely resembling one-hot vectors. This indicates that the output distributions are highly compact, making it difficult for a knockoff student to surpass LS regularization when attempting to distill knowledge from a CMIM-trained model. This increased compactness plays a crucial role in enhancing the model's resistance to knowledge distillation.

Lastly, we clarify the distinction between "highly concentrated output clusters" and "overly confident predictions". A highly concentrated output cluster does not necessarily imply that the model produces overly confident predictions. This is because the clusters can be concentrated around points that are not close to one-hot labels (the corners of the probability simplex). As a result, the model can have concentrated outputs without being overly confident. These are two separate concepts.

To illustrate this, we train an RN50 model on the CIFAR-100 dataset using three different methods: CMIM, CE, and LS. After training, we evaluate the model's confidence by measuring the average entropy of its output probability vectors on the test dataset. The results, presented in Table 4, indicate that the entropy for CMIM is higher than that for CE, demonstrating that CMIM produces less confident output probabilities.

Additionally, our experiments in Table 1 further support the above-mentioned claim: the CMIM method generates more concentrated output clusters while also improving the model's accuracy on the held-out test dataset compared to models trained with the conventional CE loss. This observation is consistent with the

 $<sup>^5\</sup>mathrm{In}$  Appendix L, we present the visualization of three different projected probability clusters for ResNet-50 trained on CIFAR-100.



Table 4: Entropy Value of the Models Trained by CMIM, CE and LS

Figure 3: The CMI values for the models trained by different KD-resistance defence methods Vs. power transform value  $\alpha$  for (a) ResNet-50, and (b) VGG16 trained on CIFAR-100 dataset. Compared to all other defense methods in the literature, CMIM effectively reduces peak CMI values under temperature scaling, which prevents the teacher model from being distilled by the knockoff student.

findings of Yang et al. (2025), which demonstrate that training DNNs to produce highly concentrated output clusters can lead to improved test accuracy.

#### 5.4 Why Prior Defense Methods Can be Made Distillable?

In this section, we address the question of why DNNs trained using all prior KD-resistance defense methods can still be made distillable, as demonstrated in Section 5.1. The key reason behind this lies in the fundamental characteristics of these defense methods and their susceptibility to distillation under certain conditions. Specifically, by appropriately tuning the power transform parameter  $\alpha$ , the models trained using these defense methods can attain a relatively high CMI value in comparison to our proposed method, CMIM (as illustrated in Figure 3). As though, compared to the CMIM model, all other defense methods provide more information for the student to leverage in improving their own performance.

This suggests that, despite their initial resistance, the defense methods fail to enforce undistillability rigorously across all distillation settings. When the correct  $\alpha$  value is selected during the distillation process, a logit-based KD approach can leverage this property to effectively distill knowledge from these supposedly resistant DNNs. Consequently, these models can still be exploited to produce distilled students that outperform the baseline LS student, demonstrating that prior defense strategies do not provide robust protection against all KD methods.

### 6 Conclusion and Future Directions

In this paper, from an information-theoretic perspective, we proposed a defence method against the threat posed by knockoff students attempting to steal the IP of pre-trained DNNs via logit-based KD methods. In particular, we proposed to minimize the CMI of the protected DNN across different power transform hyperparameter values  $\alpha$ , while minimizing the conventional CE loss simultaneously. We referred to model trained by these framework as CMIM models. By conducting a series of experiments, we showed that, unlike the prior defense methods proposed in the literature, a knockoff student cannot render CMIM models distillable. In addition, we showed that the models trained by CMIM achieve higher classification accuracy compared to those trained by CE loss. Despite these promising results, our work has certain limitations. First, the evaluation of CMIM models is primarily empirical, as providing a formal theoretical proof of undistillability remains an open challenge. Second, our approach introduces additional computational overhead compared to the conventional training using CE loss.

For future work, we aim to extend the CMIM method beyond standard classification tasks. Potential directions include adapting it for multi-label classification, regression problems, and safeguarding the intellectual property of cutting-edge models such as LLMs, CLIP, and diffusion models. Another promising direction is to enhance the model's undistillability by adapting CMIM for integration with a contrastive learning framework, which regularizes the latent space rather than the output probability space. By broadening the scope of CMIM, we hope to further enhance its applicability and effectiveness across a wider range of machine learning paradigms.

### References

- Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and selfdistillation in deep learning. arXiv preprint arXiv:2012.09816, 2020.
- Jure An, Doetsch Peter, and Pylvänäinen Thomas. Relation knowledge distillation. In International Conference on Learning Representations, 2021.
- Rohan Anil, Gabriel Pereyra, Alexandre Passos, Robert Ormandi, George E Dahl, and Geoffrey E Hinton. Large scale distributed neural network training through online distillation. *arXiv preprint* arXiv:1804.03235, 2018.
- Kenneth Borup and Lars N Andersen. Even your teacher needs guidance: Ground-truth targets dampen regularization imposed by self-distillation. Advances in Neural Information Processing Systems, 34:5316– 5327, 2021.
- Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 535–541, 2006.
- Anda Cheng and Jian Cheng. Apgp: Accuracy-preserving generative perturbation for defending against model cloning attacks. In ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 1–5, 2023. doi: 10.1109/ICASSP49357.2023.10094956.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pp. 248–255. Ieee, 2009.
- Kangning Guo, Hu Shengyuan, Yan Junjie, Liu Xin, Xu Dongbao, and Wang Ningbo. Logit-like knowledge distillation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 35, pp. 10186–10193, 2021.
- Shayan Mohajer Hamidi, Xizhen Deng, Renhao Tan, Linfeng Ye, and Ahmed Hussein Salamah. How to train the teacher model for effective knowledge distillation. In *European Conference on Computer Vision*, pp. 1–18. Springer, 2024.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 770–778, 2016. doi: 10.1109/CVPR.2016.90.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015. URL https://arxiv.org/abs/1503.02531.
- Tao Huang, Shan You, Fei Wang, Chen Qian, and Chang Xu. Knowledge distillation from a stronger teacher. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), Advances in Neural Information Processing Systems, 2022a. URL https://openreview.net/forum?id=157Usp\_kbi.

- Tao Huang, Shan You, Fei Wang, Chen Qian, and Chang Xu. Knowledge distillation from a stronger teacher. Advances in Neural Information Processing Systems, 35:33716–33727, 2022b.
- Surgan Jandial, Yash Khasbage, Arghya Pal, Vineeth N. Balasubramanian, and Balaji Krishnamurthy. Distilling the undistillable: Learning from a nasty teacher. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (eds.), Computer Vision – ECCV 2022, pp. 587–603, Cham, 2022. Springer Nature Switzerland. ISBN 978-3-031-19778-9.
- Reyhan Kevser Keser and Behcet Ugur Toreyin. Averager student: Distillation from undistillable teacher, 2023. URL https://openreview.net/forum?id=4isz71\_aZN.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). University of Toronto, 2012. URL http://www.cs.toronto.edu/~kriz/cifar.html.
- Ya Le and Xuan S. Yang. Tiny imagenet visual recognition challenge. 2015. URL https://api.semanticscholar.org/CorpusID:16664790.
- Yann LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In Neural networks: Tricks of the trade, pp. 9–50. Springer, 2002.
- Taesung Lee, Benjamin Edwards, Ian Molloy, and Dong Su. Defending against neural network model stealing attacks using deceptive perturbations. In 2019 IEEE Security and Privacy Workshops (SPW), pp. 43–49, 2019. doi: 10.1109/SPW.2019.00020.
- Haoyu Ma, Tianlong Chen, Ting-Kuei Hu, Chenyu You, Xiaohui Xie, and Zhangyang Wang. Undistillable: Making a nasty teacher that cannot teach students. In *International Conference on Learning Representa*tions, 2021. URL https://openreview.net/forum?id=0zvfm-nZqQs.
- Haoyu Ma, Yifan Huang, Tianlong Chen, Hao Tang, Chenyu You, Zhangyang Wang, and Xiaohui Xie. Stingy teacher: Sparse logits suffice to fail knowledge distillation, 2022. URL https://openreview.net/forum?id=ae7BJI0xkxH.
- Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In Proceedings of the European conference on computer vision (ECCV), pp. 116–131, 2018.
- Aditya K Menon, Ankit Singh Rawat, Sashank Reddi, Seungyeon Kim, and Sanjiv Kumar. A statistical perspective on distillation. In *International Conference on Machine Learning*, pp. 7632–7642. PMLR, 2021.
- Hossein Mobahi, Mehrdad Farajtabar, and Peter Bartlett. Self-distillation amplifies regularization in hilbert space. Advances in Neural Information Processing Systems, 33:3351–3361, 2020.
- Daniel Moldovan, Ionescu Bogdan, Drimbă Alexandru, and Marius Popescu. Path-kg: Knowledge distillation with path-level guidance. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 1709–1718, 2019.
- Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? Advances in neural information processing systems, 32, 2019.
- Tribhuvanesh Orekondy, Bernt Schiele, and Mario Fritz. Knockoff nets: Stealing functionality of black-box models. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 4954–4963, 2019.
- Tribhuvanesh Orekondy, Bernt Schiele, and Mario Fritz. Prediction poisoning: Towards defenses against dnn model stealing attacks. In *International Conference on Learning Representations*, 2020. URL https: //openreview.net/forum?id=SyevYxHtDB.

- Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3967–3976, 2019.
- Gabriel Pereyra, George Tucker, Jan Chorowski, Lukasz Kaiser, and Geoffrey Hinton. Regularizing neural networks by penalizing confident output distributions, 2017. URL https://openreview.net/forum?id=HkCjNI5ex.
- Mary Phuong and Christoph Lampert. Towards understanding knowledge distillation. In International conference on machine learning, pp. 5142–5151. PMLR, 2019.
- Radim Rehurek and Petr Sojka. Gensim–python framework for vector space modelling. NLP Centre, Faculty of Informatics, Masaryk University, Brno, Czech Republic, 3(2), 2011.
- Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. The Annals of Mathematical Statistics, 22(3):400 407, 1951. doi: 10.1214/aoms/1177729586. URL https://doi.org/10.1214/aoms/1177729586.
- Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. arXiv preprint arXiv:1412.6550, 2014.
- Ahmed H Salamah, Shayan Mohajer Hamidi, and En-Hui Yang. A coded knowledge distillation framework for image classification based on adaptive jpeg encoding. *Pattern Recognition*, 158:110966, 2025.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4510–4520, 2018.
- Reza Shokri and Vitaly Shmatikov. Privacy-preserving deep learning. In 2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton), pp. 909–910, 2015. doi: 10.1109/ALLERTON.2015.7447103.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- Zi Wang, Chengcheng Li, and Husheng Li. Adversarial training of anti-distilled neural network with semantic regulation of class confidence. In 2022 IEEE International Conference on Image Processing (ICIP), pp. 3576–3580, 2022. doi: 10.1109/ICIP46576.2022.9897169.
- En-hui Yang and Linfeng Ye. Markov knowledge distillation: Make nasty teachers trained by self-undermining knowledge distillation fully distillable. In *European Conference on Computer Vision*. Springer, 2024.
- En-hui Yang, Shayan Mohajer Hamidi, Linfeng Ye, Renhao Tan, and Beverly Yang. Conditional mutual information constrained deep learning: Framework and preliminary results. In 2024 IEEE International Symposium on Information Theory (ISIT), pp. 569–574, 2024. doi: 10.1109/ISIT57864.2024.10619241.
- En-Hui Yang, Shayan Mohajer Hamidi, Linfeng Ye, Renhao Tan, and Beverly Yang. Conditional mutual information constrained deep learning for classification. *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–13, 2025. doi: 10.1109/TNNLS.2025.3540014.
- Linfeng Ye, Shayan Mohajer Hamidi, Renhao Tan, and En-hui Yang. Bayes conditional distribution estimation for knowledge distillation based on conditional mutual information. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=yV6wwEbtkR.
- Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. Decoupled knowledge distillation. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pp. 11953–11962, 2022.
- Kaixiang Zheng and ENHUI YANG. Knowledge distillation based on tranformed teacher matching. In *The Twelfth Interna2ional Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=MJ74K7uDGG1.

### **A** Appendix

### **B** A Thorough Study of Related Works

#### B.1 Logit-based KD methods

Knowledge Distillation (KD) has become a cornerstone technique for compressing large teacher models into smaller student models. This section reviews key logit-based KD methods and explores their advancements.

Hinton et al. (2015) introduced the foundational concept of KD, using KL divergence to match the student's softmax outputs to the teacher's. This work laid the groundwork for logit-based methods, as the softmax output directly relates to the logits. Later, Moldovan et al. (2019) proposed Path-KD, a method that utilizes the paths leading to the correct class in both teacher and student models for distillation. While not directly logit-based, it demonstrates the effectiveness of aligning decision-making processes. Guo et al. (2021) proposed Logit-Like Distillation, addressing the capacity gap by matching the ranking of logits instead of their exact values [4]. This approach allows the student to learn the essential ordering of classes even with limited capacity. An et al. (2021) proposed relation knowledge distillation (RKD), focusing on aligning relationships between class logits rather than individual values. This approach improves the student's ability to generalize to unseen data. Zhao et al. (2022) introduced decoupled knowledge distillation (DKD), where it decouples the classical KD loss into two parts: target class knowledge distillation and non-target class knowledge distillation. Huang et al. (2022b) proposed DIST where they designed a KD method to distill better from a stronger teacher; indeed they claim that preserving the relations between the predictions of teacher and student would suffice for an effective KD. Borup & Andersen (2021) provided theoretical arguments for the importance of weighting the teacher outputs with the ground-truth targets when performing self-distillation with kernel ridge regressions along with a closed form solution for the optimal weighting parameter. Hamidi et al. (2024) propose to train the teacher model using mean square loss instead of the conventional cross entropy loss. Salamah et al. (2025) proposed a KD framework that leverages adaptive JPEG encoding to make the teacher's output responses less peaky for image classification tasks.

#### B.2 Defense methods against Logit-based KD

As also discussed in Section 2, the defense methods against the threat posed by knockoff students attempting to steal the IP of pre-trained DNNs via logit-based KD methods can be categorized into two approaches. Here, we elaborate on these two approaches.

(I) Model stealing resistant training: In this approach, DNNs are trained to reduce the accuracy of knockoff students while maintaining the original classification accuracy of the model. In particular, Ma et al. (2021) proposed a training algorithm named self-undermining KD to create nasty teachers (NT) that prevent knowledge leakage and unauthorized model stealing through KD, without compromising model accuracy. The nasty teacher is trained by minimizing the following objective function:

$$\mathcal{L}^{NT} = \mathsf{H}(y, q_x) - \epsilon \, \mathrm{KL}(\tilde{q}_x, q_x), \tag{30}$$

where  $\tilde{q}_x$  is a output of a pre-trained standard model.

Subsequently, Wang et al. (2022) proposed semantic nasty teachers (SNT) which improve the model stealing resistance of NT by disentangling semantic relationships in the output logits during teacher model training, which is crucial for successful KD.

(II) post-training defence methods: The aim of these approaches is to deceive the knockoff by imposing minimal perturbations to the model's predictions. Lee et al. (2019) tested a variety of possible perturbation forms, and found that the reverse sigmoid perturbation (RSP) to be the most effective one. Orekondy et al. (2020) introduced maximizing angular deviation (MAD), a technique that perturbs the output probabilities, leading to an adversarial gradient signal that deviates significantly from the original gradient of the knockoff. To this end, they applied a randomly initialized model as the surrogate for the potential knockoff. More recently, Cheng & Cheng (2023) proposed a plug-and-play generative perturbation model, dubbed as accuracy

preserving generative perturbation (APGP), which can effectively defend KD-based model cloning, while preserve the model utility.

#### B.3 Attack Methods Using Logit-based KD

Jandial et al. (2022) sought to circumvent the defense of nasty teachers and steal (or extract) its information. Specifically, they analyzed nasty teacher from two different angles and subsequently leverage them carefully to develop simple yet efficient methodologies, named as HTC and SCM, which enhance learning from nasty teacher.

In AVG (Keser & Toreyin, 2023), the authors noted that undistillable teachers exhibit multiple peaks in their softmax response, which are transferred to the student models. These peaks are considered to be the primary factor that misleads the student models. To mitigate the influence of the multiple peaks in the softmax response of teachers, they proposed transferring the mean of features with the same labels as the soft labels.

Orekondy et al. (2019) introduced a technique called "Knockoff Nets" that allows an attacker to steal the functionality of black-box models. Remarkably, the attacker only needs to interact with the model by feeding it input data and observing the resulting predictions. By training a new model ("knockoff") on these input-prediction pairs, the attacker can create a copycat model that performs similarly to the original black box.

### C Why LS reduce the knockoff student's accuracy?

Original aiming to prevent overfitting and improve generalization, label smoothing was observed by Müller et al. (2019) to reduce the accuracy of the knockoff student. The researchers found that "label smoothing encourages examples to lie in tight, equally separated clusters". Consequently, label smoothing reduces the contextual information in the teacher model's output (Yang et al., 2024).

### D Power Transformation of Mutual Information

The following theorem implies that for each label y,  $I(X; \hat{Y}^{\alpha}|Y = y)$  as a function of  $\alpha$  is continuously differentiable.

**Theorem 2.** Let (X, Z) be a pair of random variables, where Z is discrete, and X can be either discrete or continuous. Let  $P_{Z|X}[\cdot|x]$  denote the conditional probability distribution of Z given X = x. Additionally, let  $P_{Z|X}^{\alpha}[\cdot|x]$  denote the power transformed version of  $P_{Z|X}[\cdot|x]$  with power  $\alpha$ ,  $Z^{\alpha}$  denote the random variable the conditional distribution of which given X = x is  $P_{Z|X}^{\alpha}[\cdot|x]$ , and  $q^{\alpha}$  denote the probability distribution of  $Z^{\alpha}$ . Then, the following holds

$$\frac{\partial \operatorname{I}(X;Z^{\alpha})}{\partial \alpha} = \frac{1}{\alpha} \sum_{x} P_{X}[x] \Big\{ \Big( m_{2}(P_{Z|X}^{\alpha}[\cdot|x]) - m_{1}^{2}(P_{Z|X}^{\alpha}[\cdot|x]) \Big) - \operatorname{Cov}(P_{Z|X}^{\alpha}[\cdot|x],q^{\alpha}) \Big\},\tag{31}$$

where for probability vectors P and Q,

$$m_1(P) \stackrel{\Delta}{=} \sum_j P[j] \big( -\ln(P[j]) \big) = \mathsf{H}(P), \quad \text{(Shannon entropy)} \tag{32a}$$

$$m_2(P) \stackrel{\Delta}{=} \sum_j P[j] (-\ln(P[j]))^2$$
, (Second moment) (32b)

$$\operatorname{Cov}(P,Q) \stackrel{\Delta}{=} \sum_{j} P[j] \Big( -\ln(P[j]) - m_1(P) \Big) \Big( -\ln(Q[j]) - \sum_{i} P[i](-\ln(Q[i])) \Big).$$
(32c)

*Proof.* To simplify our notation, we denote the conditional distributions  $P_{Z|X}[\cdot|x]$  and  $P_{Z|X}^{\alpha}[j|x]$  by  $p_x$  and  $p_x^{\alpha}$ , respectively. Decompose  $I(X; Z^{\alpha})$  as follows

$$I(X; Z^{\alpha}) = H(Z^{\alpha}) - H(Z^{\alpha}|X)$$
  
= H(q^{\alpha}) -  $\sum_{x} P[x]H(p_{x}^{\alpha})$  (33)

where for any random variables U and V, H(V) and H(V|U) are the entropy of V and the conditional entropy of V given U, respectively, and H(p) denotes the entropy of the probability distribution p. Then the partial derivative of  $I(X; Z^{\alpha})$  w.r.t.  $\alpha$  is equal to

$$\frac{\partial \mathbf{I}(X;Z^{\alpha})}{\partial \alpha} = \frac{\partial \mathbf{H}(q^{\alpha})}{\partial \alpha} - \sum_{x} P[x] \frac{\partial \mathbf{H}(p_{x}^{\alpha})}{\partial \alpha}.$$
(34)

To continue, we first compute the partial derivative in the second term of the RHS of equation 34

$$\frac{\partial \mathrm{H}(p_{x}^{\alpha})}{\partial \alpha} = \frac{-\partial \sum_{j} p_{x}^{\alpha}[j] \ln(p_{x}^{\alpha}[j])}{\partial \alpha} = -\sum_{j} \left( \ln(p_{x}^{\alpha}[j]) + 1 \right) \frac{\partial p_{x}^{\alpha}[j]}{\partial \alpha} \\
= -\sum_{j} \left( \ln(p_{x}^{\alpha}[j]) + 1 \right) \\
\times \frac{(p_{x}[j])^{\alpha} \ln(p_{x}[j]) \left( \sum_{i} (p_{x}[i])^{\alpha} \right) - (p_{x}[j])^{\alpha} \left( \sum_{i} (p_{x}[i])^{\alpha} \ln p_{x}[i] \right)}{\left( \sum_{i} (p_{x}[i])^{\alpha} \right)^{2}} \\
= -\sum_{j} \left( \ln(p_{x}^{\alpha}[j]) + 1 \right) \left( p_{x}^{\alpha}[j] \left( \ln(p_{x}[j]) - \sum_{i} p_{x}^{\alpha}[i] \ln(p_{x}[i]) \right) \right) \\
= \frac{-1}{\alpha} \sum_{j} \left( \ln(p_{x}^{\alpha}[j]) + 1 \right) p_{x}^{\alpha}[j] \left( \ln(p_{x}[j])^{\alpha} - \sum_{i} p_{x}^{\alpha}[i] \ln(p_{x}^{\alpha}[i]) \right) \\
= \frac{-1}{\alpha} \left( \sum_{j} p_{x}^{\alpha}[j] \left( \ln(p_{x}^{\alpha}[j]) \right)^{2} - \left( \sum_{j} p_{x}^{\alpha}[j] \ln(p_{x}^{\alpha}[j]) \right) \left( \sum_{i} p_{x}^{\alpha}[i] \ln(p_{x}^{\alpha}[i]) \right) \right) \\
= \frac{-1}{\alpha} \left( m_{2}(p_{x}^{\alpha}) - m_{1}^{2}(p_{x}^{\alpha}) \right)$$
(36)

Note that

$$q^{\alpha} = \sum_{x} P[x] p_x^{\alpha}.$$

Then we have

where equation 37 is due to equation 35.

From Equations (36) and (38), Theorem 2 follows.

**E Proof of Theorem 1** 

Theorem 1 follows from Theorem 2 and the following lemma. Lemma 1. Let g(t) be a continuously differentiable function over  $[0, \beta]$ . Then the following holds:

$$\max_{t} g(t) = \lim_{\omega \to \infty} \frac{1}{\omega} \ln \frac{1}{\beta} \int_{0}^{\beta} \exp\left\{\omega g(t)\right\} dt.$$
(39)

*Proof.* Let  $t^*$  be an optimal point at which

$$g(t^*) = \max_t g(t).$$

For any  $\epsilon > 0$ , let  $\mathcal{N}(t^*, \epsilon)$  denote a closed interval containing  $t^*$  with length  $\epsilon$ . It is easy to verify that

$$\frac{1}{\omega} \ln \frac{1}{\beta} \int_0^\beta \exp\left\{\omega g(t)\right\} dt \le g(t^*)$$

which implies that

$$\limsup_{\omega \to \infty} \frac{1}{\omega} \ln \frac{1}{\beta} \int_0^\beta \exp\left\{\omega g(t)\right\} dt \le g(t^*).$$
(40)

On the other hand,

$$\frac{1}{\omega} \ln \frac{1}{\beta} \int_{0}^{\beta} \exp\left\{\omega g(t)\right\} dt \geq \frac{1}{\omega} \ln \frac{1}{\beta} \int_{\mathcal{N}(t^{*},\epsilon)} \exp\left\{\omega g(t)\right\} dt$$

$$\geq \frac{1}{\omega} \ln \frac{\epsilon}{\beta} \exp\left\{\omega \min_{t \in \mathcal{N}(t^{*},\epsilon)} g(t)\right\}$$

$$= \min_{t \in \mathcal{N}(t^{*},\epsilon)} g(t) + \frac{1}{\omega} \ln \frac{\epsilon}{\beta}.$$
(41)

Letting  $\omega \to \infty$  in equation 41 yields

$$\liminf_{\omega \to \infty} \frac{1}{\omega} \ln \frac{1}{\beta} \int_0^\beta \exp\left\{\omega g(t)\right\} dt \ge \min_{t \in \mathcal{N}(t^*, \epsilon)} g(t).$$
(42)

Note that equation 42 is valid for any  $\epsilon > 0$ . Letting  $\epsilon \to 0$  in equation 42, we have

$$\liminf_{\omega \to \infty} \frac{1}{\omega} \ln \frac{1}{\beta} \int_0^\beta \exp\left\{\omega g(t)\right\} dt \ge g(t^*).$$
(43)

Then equation 39 follows from equation 40 and equation 43. This completes the proof of Lemma 1.  $\Box$ 

### F Datasets description

- CIFAR-100 (Krizhevsky et al., 2012) dataset contains 50K training and 10K test color images, each with size 32 × 32, categorized into 100 classes.
- TinyImageNet (Le & Yang, 2015) contains 120K color images across 200 classes, each with a resolution of 64 × 64 pixels. For each class, there are 500 training images, 50 validation images and 50 test images.
- ImageNet (Deng et al., 2009) is a large-scale dataset used in visual recognition tasks, containing around 1.2 million training and 50K validation images.

### G Experiments setup

All experiments detailed in this paper were conducted using a publicly available national high-performance computer. For each experiment, we utilized 16 CPU cores, 64 GB of memory, and one NVIDIA V100 GPU. The software environment comprised Python 3.10, PyTorch 1.13, and CUDA 11.

For all experiments, including defenses and attacks, the SGD optimizer (Robbins & Monro, 1951; LeCun et al., 2002) with a learning rate of 0.1 is used unless otherwise specified.

For the CIFAR-100 and TinyImageNet datasets, we train the model for 200 epochs, decaying the learning rate by 0.1 at epochs 60, 120, 160.

For ImageNet, we follow the standard PyTorch practice  $^{6}$ .

The batch size is 128 for both CIFAR-100 and TinyImageNet, and 256 for ImageNet.

To get the accuracy that a knockoff student can achieve using label smoothing, we have tested a wide spectrum of label smoothing factor  $\epsilon = \{0.01, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\}$ , and selected the value that yielded a classification accuracy exceeding that of all knockoff students.

In the CMIM method, we set T = 20 and tested  $\lambda = 0.1, 0.25, 0.5, 1$ , selecting the value that minimized the CMI value while maintaining or improving classification accuracy.

#### G.1 Defense setup

We used the following parameters and settings for the defense models used in Section 5.

#### G.1.1 Defense setup on CIFAR-100 and TinyImagenet

- 1. **MAD**: We employ a randomly initialized VGG-8 as adversary's architecture, and following the implementation of MAD-argmax.
- 2. **APGP**: We apply a 3 layer MLP with residual connection as the generative model and set  $\lambda = 0.1$  for all experiments.

<sup>&</sup>lt;sup>6</sup>https://github.com/pytorch/vision/tree/main/references/classification

- 3. **RSP**: We use  $\alpha = 1$  and  $\beta = 20$  for all the experiments.
- 4. NT: To ensure a acceptable accuracy sacrifice, we test three different  $\epsilon$  values and select the largest one that results in an accuracy drop of less than 0.5%. Specifically, we use  $\epsilon = 0.01$  for ResNet-50 and  $\epsilon = 0.005$  for VGG-16 on the CIFAR-100 dataset, while for TinyImageNet, we use  $\epsilon = 0.001$  for both ResNet-34 and ResNet-50.
- 5. **SNT**: We use the pretrained word2vect model namely "fasttext-wiki-news-subwords-300" provided by Gensim (Rehurek & Sojka, 2011), and set  $\lambda = 0.2$  for all experiments.
- 6. **ST**: We use the teacher model trained by CE as the underlying model, and use the sparse ratio of 10% as suggested in their paper for all experiments.
- 7. LS: We apply label smoothing factor 0f 0.05 for all experiments.
- 8. **ISTM**: We set the binary search parameters to  $T_b = 20$  and  $\alpha_{\text{max}} = 2000$ . We use  $\lambda = 0.2$  for ResNet-50 and  $\lambda = 0.5$  for VGG-16 on the CIFAR-100 dataset, while for TinyImageNet, we use  $\lambda = 0.1$  for ResNet-34 and  $\lambda = 0.5$  ResNet-50.

#### G.1.2 Defense setup on Imagenet

- 1. **ST**: We use the teacher model trained by CE as the underlying model, and use the sparse ratio of 10% as suggested in their paper for all experiments.
- 2. **ISTM**: We set the binary search parameters to  $T_b = 20$  and  $\alpha_{\text{max}} = 2000$ . We use  $\lambda = 0.2$  for all the experiments.

#### G.2 Attack setup

#### G.2.1 Attack setup on CIFAR-100 and TinyImagenet

We use power transform parameter  $\alpha = 0.25$  (or equivalently T = 4) for all experiments unless otherwise specified.

- 1. **KD**: We set the CE-KL trade-off coefficient to  $\lambda = 0.9$ .
- 2. **MKD**: We use the intrinsic dimension of 3 for CIFAR-100, and 5 for TinyImageNet. We employed the Adam optimizer (Kingma & Ba, 2014) with learning rate  $10^{-3}$  for the trainable Markov transform.
- 3. **DKD**: We test alpha, beta pairs of  $\{1, 4\}$  and  $\{2, 8\}$ , and report the one with best accuracy.
- 4. **DIST**: We  $use\beta = 1.0, \gamma = 1.0, \tau = 1.0$  for all experiments.
- 5. **HTC**: We use  $\alpha = 0.05(T = 20)$ ,  $\lambda = 0.01$  for all experiments.
- 6. **AVG**:  $\lambda = 0.9$ .
- 7. Knockoff: We follow the implementation of the original paper.

#### G.2.2 Attack setup on Imagenet

We use  $\alpha = 1$  (T = 1) for all experiments unless otherwise specified.

- 1. **KD**:  $\lambda = 0.9$ .
- 2. **MKD**: We use the intrinsic dimension of 16. We employed the Adam optimizer (Kingma & Ba, 2014) with learning rate  $10^{-3}$  for the trainable Markov transform.
- 3. **DKD**: We test alpha, beta pairs of  $\{1, 4\}$  and  $\{2, 8\}$ , and report the one with best accuracy.

### H Accuracy of Protected models

In this section, we report the top-1 accuracy of some additional models that are trained using CMIM and compare them with those trained by CE method. To this end, we use 10 well-known models for CIFAR-100 dataset namely ResNet (RN)-{18, 34, 50, 101, 152}, SqueezeNet (SQN), ResNext (RNXT) 50, MobileNet (MN), Xception (XCP), DenseNet (DN) 121; and 2 models namely RN-{34, 50} for TinyImageNet and ImageNet. We follow the same training recipe as the one in Section 5.1. The results for CIFAR-100 and (Tiny-)ImageNet are listed in Table 5 and Table 6, respectively. As seen, the top-1 accuracy for all models trained by CMIM is consistently higher than those trained by CE counterpart, with the gain up to 1.15%.

Table 5: Top-1 accuracy (%) of models trained by CE and CMIM methods on CIFAR-100.

CIFAR-100										
Model	CE	CMIM	Model	CE	CMIM					
RN18	76.05	77.20	SQN	69.32	70.64					
RN34	77.20	77.54	RNXT50	78.71	79.12					
RN50	77.81	77.93	MN	67.26	67.51					
RN101	79.07	79.12	XCP	77.37	77.64					
RN152	79.21	79.43	DN121	79.16	79.33					

Table 6: Top-1 accuracy (%) of models trained by CE and CMIM methods on TinyImageNet and ImageNet.

Tir	nyImage	Net	ImageNet				
Model	CE	CMIM	Model	CE	CMIM		
RN34	65.39	65.99	RN34	73.31	73.69		
RN50	66.14	66.93	RN50	76.15	76.40		

## I Variance of Table 1

Table 7: Top-1 accuracy (%) and variance of the knockoff student on CIFAR-100 and TinyImageNet dataset (averaged over 3 runs)

	CIFAR-100										
Defense	Model	K-student	LS	KD	MKD	DKD	DIST	HTC	AVG	Knockoff	
	VOOIC	VGG11	$71.94{\pm}0.09$	$68.55 \pm 0.20$	$72.08 \pm 0.29$	$53.32 {\pm} 0.38$	$69.21 \pm 0.09$	$71.19 {\pm} 0.03$	$70.03 {\pm} 0.07$	$61.44 \pm 0.06$	
MAD	VGG10	SNV2	$72.65 {\pm} 0.18$	$72.50 {\pm} 0.11$	$72.46 {\pm} 0.13$	$7.64 \pm 0.70$	$69.91\pm0.18$	$71.37 {\pm} 0.24$	$72.86{\pm}0.18$	$70.87 {\pm} 0.26$	
MAD	DNFO	VGG11	$71.94{\pm}0.09$	$72.00 \pm 0.21$	$72.04 \pm 0.13$	$54.29 \pm 0.52$	$71.57 \pm 0.31$	$70.76 {\pm} 0.20$	$70.73 \pm 0.22$	$61.73 \pm 0.23$	
	ILIN50	RN18	$78.76 {\pm} 0.08$	$77.76 {\pm} 0.23$	$78.79 {\pm} 0.23$	$43.73 {\pm} 0.55$	$73.76 {\pm} 0.10$	$77.89 {\pm} 0.25$	$78.61 {\pm} 0.15$	$73.92{\pm}0.19$	
VCC	VCC16	VGG11	$71.94{\pm}0.09$	$71.92 {\pm} 0.19$	$72.27 \pm 0.21$	$27.24{\pm}0.54$	$69.25 \pm 0.14$	$70.08 {\pm} 0.23$	$72.01 \pm 0.20$	$45.98 \pm 0.45$	
APCP	VGG10	SNV2	$72.65 \pm 0.18$	$73.10 {\pm} 0.27$	$73.75 {\pm} 0.23$	$12.52{\pm}0.30$	$71.04{\pm}0.31$	$71.66 {\pm} 0.13$	$73.20{\pm}0.27$	$9.48 \pm 0.73$	
AIGI	RN50	VGG11	$71.94{\pm}0.09$	$71.91 \pm 0.17$	$72.11 \pm 0.23$	$9.74 \pm 0.86$	$69.48 \pm 0.11$	$71.38 \pm 0.25$	$71.92 \pm 0.14$	$34.71 \pm 0.30$	
	10100	RN18	$78.76 \pm 0.08$	$78.04 \pm 0.21$	$79.06 \pm 0.14$	$62.71 \pm 0.29$	$77.32 \pm 0.13$	$77.82 \pm 0.13$	$77.90 \pm 0.15$	$2.57 \pm 0.95$	
	VGG16	VGG11	$71.94 \pm 0.09$	$71.42 \pm 0.24$	$72.04 \pm 0.13$	$70.22 \pm 0.19$	$70.80 \pm 0.17$	$70.40 \pm 0.17$	$71.56 \pm 0.06$	$31.04 \pm 0.62$	
RSP	Vadio	SNV2	$72.65 \pm 0.18$	$73.55 \pm 0.34$	$72.95 \pm 0.36$	$67.45 \pm 0.20$	$72.19 \pm 0.44$	$71.46 \pm 0.41$	$72.27 \pm 0.27$	$26.09 \pm 0.40$	
1651	RN50	VGG11	$71.94 \pm 0.09$	$71.97 \pm 0.17$	$72.01 \pm 0.20$	$69.53 \pm 0.12$	$72.18 \pm 0.21$	$70.87 \pm 0.14$	$70.85 \pm 0.17$	$46.68 \pm 0.60$	
	101000	RN18	$78.76 \pm 0.08$	$77.78 \pm 0.09$	$77.79 \pm 0.16$	$77.01 \pm 0.09$	$78.88 \pm 0.21$	$78.00 \pm 0.26$	$78.13 \pm 0.12$	$55.86 \pm 0.18$	
	VGG16	VGG11	$71.94 \pm 0.09$	$71.40\pm0.34$	$73.44 \pm 0.16$	$71.47 \pm 0.14$	$71.33\pm0.18$	$70.77 \pm 0.23$	$71.58 \pm 0.09$	$63.56 \pm 0.16$	
NT	10010	SNV2	$72.65 \pm 0.18$	$72.44 \pm 0.43$	$72.70 \pm 0.35$	$6.24 \pm 0.51$	$72.04 \pm 0.19$	$70.75 \pm 0.13$	$72.83 \pm 0.20$	$6.32 \pm 0.26$	
	RN50	VGG11	$71.94 \pm 0.09$	$72.01 \pm 0.25$	$72.03 \pm 0.19$	$71.55 \pm 0.36$	$71.88 \pm 0.31$	$70.16 \pm 0.29$	$71.94 \pm 0.18$	$62.94 \pm 0.24$	
	10100	RN18	$78.76 \pm 0.08$	$78.41 \pm 0.25$	$78.92 \pm 0.14$	$79.26 \pm 0.29$	$78.99 \pm 0.14$	$77.94 \pm 0.22$	$78.33 \pm 0.05$	$68.96 \pm 0.18$	
	VGG16	VGG11	$71.94 \pm 0.09$	$72.06 \pm 0.22$	$72.28 \pm 0.12$	$4.92 \pm 0.22$	$71.98 \pm 0.18$	$70.60 \pm 0.13$	$71.63 \pm 0.10$	$64.08 \pm 0.19$	
SNT	10010	SNV2	$72.65 \pm 0.18$	$72.94 \pm 0.41$	$73.17 \pm 0.13$	$72.78 \pm 0.20$	$72.22 \pm 0.24$	$71.22 \pm 0.18$	$72.74 \pm 0.20$	$6.22 \pm 0.59$	
B	RN50	VGG11	$71.94 \pm 0.09$	$72.02 \pm 0.19$	$72.12 \pm 0.39$	$72.32 \pm 0.33$	$71.70 \pm 0.39$	$70.66 \pm 0.17$	$71.65 \pm 0.20$	$62.94 \pm 0.29$	
	101100	RN18	$78.76 \pm 0.08$	$78.25 \pm 0.05$	$78.48 \pm 0.24$	$78.82 \pm 0.30$	$78.14 \pm 0.28$	$78.45 \pm 0.15$	$78.38 \pm 0.13$	$67.71 \pm 0.20$	
	VGG16	VGG11	$71.94 \pm 0.09$	$72.09 \pm 0.23$	$72.01 \pm 0.14$	$71.63 \pm 0.16$	$71.93 \pm 0.12$	$71.16 \pm 0.27$	$71.63 \pm 0.18$	$63.32 \pm 0.14$	
ST	10010	SNV2	$72.65 \pm 0.18$	$72.64 \pm 0.15$	$72.67 \pm 0.21$	$70.53 \pm 0.48$	$72.24 \pm 0.39$	$71.32 \pm 0.38$	$72.42 \pm 0.12$	$69.46 \pm 0.28$	
~ -	RN50	VGG11	$71.94 \pm 0.09$	$72.00 \pm 0.19$	$72.13 \pm 0.13$	$71.62 \pm 0.24$	$71.76 \pm 0.29$	$70.54 \pm 0.33$	$71.73 \pm 0.11$	$65.43 \pm 0.24$	
		RN18	$78.76 \pm 0.08$	$78.96 \pm 0.26$	$79.02 \pm 0.06$	$78.35 \pm 0.09$	$78.31 \pm 0.14$	$78.36 \pm 0.25$	$78.81 \pm 0.14$	$72.87 \pm 0.08$	
	VGG16	VGG11	$71.94 \pm 0.09$	$71.90 \pm 0.18$	$72.00 \pm 0.06$	$71.57 \pm 0.26$	$70.89 \pm 0.12$	$70.66 \pm 0.17$	$71.76 \pm 0.10$	$63.49 \pm 0.21$	
LS		SNV2	$72.65 \pm 0.18$	$72.87 \pm 0.28$	$73.52 \pm 0.25$	$70.01 \pm 0.32$	$71.49 \pm 0.38$	$71.70 \pm 0.42$	$73.01 \pm 0.27$	$65.20 \pm 0.14$	
	RN50	VGG11	$71.94 \pm 0.09$	$71.82 \pm 0.28$	$71.99 \pm 0.16$	$71.95 \pm 0.33$	$70.77 \pm 0.39$	$70.86 \pm 0.24$	$71.88 \pm 0.16$	$62.29 \pm 0.10$	
		RN18	$78.76 \pm 0.08$	$77.72 \pm 0.30$	$77.82 \pm 0.12$	$79.37 \pm 0.19$	$78.33 \pm 0.06$	$78.31 \pm 0.21$	$77.91 \pm 0.07$	$63.36 \pm 0.17$	
	VGG16	VGG11	$71.94 \pm 0.09$	$71.87 \pm 0.24$	$71.64 \pm 0.07$	$71.56 \pm 0.03$	$70.34 \pm 0.09$	$71.71 \pm 0.14$	$71.42 \pm 0.05$	$66.89 \pm 0.11$	
CMIM		SNV2	$72.65 \pm 0.18$	$72.53 \pm 0.21$	$71.44 \pm 0.16$	$72.46 \pm 0.20$	$71.45 \pm 0.31$	$71.59 \pm 0.24$	$71.94 \pm 0.20$	$64.45 \pm 0.24$	
	RN50	VGGII	$71.94 \pm 0.09$	$71.54 \pm 0.30$	$71.34 \pm 0.16$	$71.77 \pm 0.06$	$71.86 \pm 0.28$	$69.32 \pm 0.09$	$71.70 \pm 0.22$	$60.58 \pm 0.17$	
		RN18	$78.76 \pm 0.08$	$78.21 \pm 0.13$	$78.16 \pm 0.09$	$78.13 \pm 0.06$	$77.56 \pm 0.06$	$77.23 \pm 0.09$	$78.64 \pm 0.06$	$65.88 \pm 0.09$	
					TinyIma	igeNet					
DCD	RN34	RN18	$63.56 \pm 0.06$	$63.54 \pm 0.09$	$64.32 \pm 0.07$	$64.01 \pm 0.07$	$63.27 \pm 0.16$	$63.54 \pm 0.07$	$62.15 \pm 0.14$	$55.43 \pm 0.12$	
RSP	RN50	SNV2	$60.61 \pm 0.15$	$60.18 \pm 0.26$	$60.76 \pm 0.16$	$56.26 \pm 0.16$	$56.43 \pm 0.11$	$60.96 \pm 0.22$	$60.15 \pm 0.20$	$54.01 \pm 0.22$	
0TT	RN34	RN18	$63.56 \pm 0.06$	$63.96 \pm 0.13$	$64.12 \pm 0.07$	$63.25 \pm 0.10$	$63.51 \pm 0.14$	$63.49 \pm 0.19$	$63.84 \pm 0.08$	$57.42 \pm 0.08$	
51	RN50	SNV2	$60.61 \pm 0.15$	$61.23 \pm 0.24$	$61.36 \pm 0.14$	$60.43 \pm 0.14$	$60.32 \pm 0.24$	$60.22 \pm 0.17$	$61.13 \pm 0.13$	$55.84 \pm 0.11$	
NT	RN34	RN18	$63.56 \pm 0.06$	$63.27 \pm 0.14$	$64.49 \pm 0.17$	$64.67 \pm 0.16$	$63.43 \pm 0.20$	$63.50 \pm 0.10$	$64.43 \pm 0.11$	$53.11 \pm 0.06$	
INI	RN50	SNV2	$60.61 \pm 0.15$	$59.57 \pm 0.22$	$61.55 \pm 0.12$	$31.55 \pm 0.28$	$60.03\pm0.23$	$60.98 \pm 0.27$	$60.31\pm0.17$	$50.94 \pm 0.15$	
TC	RN34	RN18	$63.56 \pm 0.06$	$63.74 \pm 0.08$	$64.01 \pm 0.14$	$64.23 \pm 0.11$	$63.51 \pm 0.07$	$64.20\pm0.16$	$63.04\pm0.13$	$57.43 \pm 0.10$	
LS	RN50	SNV2	$60.61 \pm 0.15$	$60.32\pm0.24$	$60.93 \pm 0.15$	$60.74\pm0.26$	$60.11\pm0.28$	$60.46\pm0.14$	$60.14\pm0.21$	$52.96 \pm 0.23$	
CMIN	RN34	RN18	$63.53 \pm 0.06$	$62.89 \pm 0.03$	$63.15 \pm 0.08$	$62.94\pm0.03$	$63.28 \pm 0.05$	$61.57 \pm 0.06$	$62.96 \pm 0.06$	$56.13 \pm 0.04$	
OMIM	RN50	SNV2	$60.61 \pm 0.15$	$57.57 \pm 0.20$	$59.32 \pm 0.17$	$60.58 \pm 0.12$	$59.41 \pm 0.09$	$59.33 \pm 0.10$	$60.42\pm0.04$	$56.91 \pm 0.05$	

Upon reviewing the variance estimates, we see that certain cases, such as the (RN50, VGG11) pair on CIFAR-100, might suggest that the CMIM-trained model could be rendered distillable under naive statistical interpretations. For example, when the DIST method is applied to the CMIM model, the accuracy is reported as, which could potentially exceed the LS student accuracy of if variance is simply added to the mean.

However, this approach of directly comparing mean values with added variances does not provide an accurate or fair assessment of undistillability. To address this, we conducted a more rigorous analysis where, across five different seeds, we compared the accuracy of the knock-off Student (VGG11 in this case) trained via label smoothing and the DIST method applied to RN50 trained with CMIM.

The results of this comprehensive comparison are summarized in the following table, which demonstrates that the CMIM model achieves undistillability across all different seeds.

### J Computational Overhead

In this section we compare the computational overhead of our method with the CE counterpart on CIFAR-100 dataset.

Table 5: Variance analysis of Table 7										
	Seed 1	Seed 2	Seed 3	Seed 4	Seed 5					
LS CMIM	72.28 72.01	71.83 71.74	72.06 71.93	72.25 72.12	71.53 71.27					

Table 9: Training times of CMIM and CE for ResNet-10 and VGG-16 on the CIFAR-100 dataset.

	CE	CMIM			
RN50	4 hours 43 minutes	5 hours 13 minutes			
VGG16	2 hours 57 minutes	3 hours 25 minutes			

Note that the training time for CMIM is slightly higher than that of conventional CE method. This is primarily due to the additional inference samples required to estimate the CMI per class. We believe this is a reasonable trade-off given the significant benefits of the method. In addition, note that the number of samples N does not have any effects on the training time CMIM; this is because the power transform applied to the teacher's output probabilities, and when calculating the gradients during the backpropagation, different values of  $\alpha$  does not change the gradients.

Additionally, we note that among the existing KD protection methods in the literature, only CMIM (our method) and ST are scalable to larger datasets like ImageNet. This scalability is due to the significant computational complexity of other benchmark methods, which limits their applicability to smaller datasets.

#### Κ Ablation on Hyper-parameters

In this section, we perform an ablation study to analyze the impact of three key hyper-parameters of CMIM:  $\beta$ , the number of samples N, and  $\omega$ .

For all experiments in this study, we employ the VGG-16 and SNV2 models as the teacher-student pair, using the CIFAR-100 dataset as the evaluation benchmark.

#### K.1 Range of $\beta$

In this section, we examine the impact of  $\beta$  on the performance of CMIM. For this analysis, we fix N = 50and  $\omega = 25$ , while varying the value of  $\beta$ . The results are shown in Figure 4 show that the accuracy of the knockoff student decreases significantly when  $\beta \geq 1$ , highlighting the importance of choosing an appropriate  $\beta$  to effectively balance cross-entropy and CMI objectives.



Figure 4: The student accuracy when distilled from the teacher model trained by different  $\beta$  values.

#### K.2 Number of power samples N

In this section, we analyze the influence of the number of samples N on the performance of the knockoff student. For this study, we set  $\beta = 2$  and  $\omega = 25$ , varying N to observe its impact. The results, presented in Figure 5, reveal that the accuracy of the knockoff student declines monotonically as N increases, suggesting that a moderate sampling size is enable to capture sufficient diversity for robust estimation.



Figure 5: The student accuracy when distilled from the teacher model trained by different number of sample N.

#### K.3 Power Coefficient $\omega$

In this section, we investigate the effect of the power coefficient  $\omega$  on the knockoff student's performance. For this analysis, we fix  $\beta = 2$  and N = 25, while varying  $\omega$ . The results are summarized in Table 10. Notably, when  $\omega > 30$ , the simulation frequently results in NaN values due to excessively large exponent values. Furthermore, the table shows that at  $\omega = 20$  or  $\omega = 30$ , the knockoff student's accuracy reaches its minimum, indicating that this value effectively approximates the behavior of  $\omega = \infty$ .

Table 10: The student accuracy when distilled from the teacher model trained by different values of  $\omega$ .

Value of $\omega$	1	2	5	10	15	20	25	30	40	50	100	200
Knock-off Student Accuracy	72.65	72.67	72.55	72.56	72.55	72.52	72.53	72.52	NaN	NaN	NaN	NaN

### L Visualization of Different Three Projected Probability Clusters



Figure 6: Visualization of three projected probability clusters for ResNet-50 trained on CIFAR-100 using (a) LS, (b) NT, and (c) CMIM.

In this section, we randomly selected three additional output clusters from the CIFAR-100 dataset, oak, raccoon, and ray, and visualized them using the method introduced in Yang et al. (2025). The same conclusion continues to hold.

### M Effect of CMIM on Model Calibration

In this section, we investigate the impact of the CMIM method on model calibration. To do so, we report the Expected Calibration Error (ECE) for ResNet-50 trained with CE, LS, NT and CMIM on the CIFAR-100 dataset. As shown in Table 11, the CMIM-trained model achieves calibration on par with its CE-trained counterpart and consistently outperforms both the LS and NT counterparts.

Table 11: Expected calibration error of three projected probability clusters for ResNet-50 trained on CIFAR-100 using CE, LS, NT and CMIM.

Defense	CE	LS	NT	CMIM
ECE	9.42%	21.48%	16.47%	11.24%
Accuracy	77.81%	78.45%	77.31%	78.72%