ENRICHING KNOWLEDGE DISTILLATION WITH INTRA-CLASS CONTRASTIVE LEARNING

Anonymous authorsPaper under double-blind review

ABSTRACT

Since the advent of knowledge distillation, much research has focused on how the soft labels generated by the teacher model can be utilized effectively. Allen-Zhu & Li (2020) points out that the implicit knowledge within soft labels originates from the multi-view structure present in the data. Feature variations within samples of the same class allow the student model to generalize better by learning diverse representations. However, in existing distillation methods, teacher models predominantly adhere to ground-truth labels as targets, without considering the diverse representations within the same class. Therefore, we propose incorporating an intra-class contrastive loss during teacher training to enrich the intra-class information contained in soft labels. In practice, we find that intra-class loss causes instability in training and slows convergence. To mitigate these issues, margin loss is integrated into intra-class contrastive learning to improve the training stability and convergence speed. Simultaneously, we theoretically analyze the impact of this loss on the intra-class distances and inter-class distances. It has been proved that the intra-class contrastive loss can enrich the intra-class diversity. Experimental results demonstrate the effectiveness of the proposed method.

1 Introduction

Knowledge distillation (KD) is a technique in deep learning where a smaller model is trained to replicate the behavior of a larger, more complex model Hinton et al. (2015); Ji et al. (2023). This approach is valuable for compressing large models Kim et al. (2018); Jin et al. (2021); Gu et al. (2023), transfer learning Vapnik et al. (2015); Noroozi et al. (2018) and enhancing the performance of smaller modelsBuciluă et al. (2006); Romero et al. (2014). Therefore, KD has gained popularity across various tasks, such as image classification Mobahi et al. (2020), natural language processing Rashid et al. (2020); Yang et al. (2020), and multimodal learning Dai et al. (2022).

Soft labels play a crucial role in knowledge distillation by offering more comprehensive information regarding the data distribution Menon et al. (2021); Zhou & Song (2021). They encapsulate the knowledge of the teacher model, encompassing not only the relative probabilities of different classes but also the intra-class variances, which are not reflected in the hard labels. Allen-Zhu & Li (2020) proposes the multi-view data assumption, which is validated by real-world datasets, and demonstrates that students can learn features of other classes present in the soft labels. These points also constitute the core hypothesis of this paper, that soft labels express the variances among samples of the same class, which is absent in ground-truth labels.

In fact, existing methods for training teacher models can be roughly divided into two categories. One targets the ground-truth (possibly with regularization) Tian et al. (2019); Chen et al. (2021a), and the other involves joint training of teacher and student models Chen et al. (2020); Xu et al.; Neitz et al. (2020); Yuan et al. (2021), where the teacher's soft labels could be adjusted based on feedback from the student. These two methods are neither designed for intra-class diversity nor provide parameters to control intra-class variation. The soft labels generated by the teacher model may fail to capture a significant amount of valuable intra-class information.

Therefore, we propose incorporating an intra-class contrastive loss as an auxiliary loss during teacher training. This can enrich the intra-class information within soft labels, preventing the soft labels from being overly similar to the ground truth due to the model's strong fitting capability. Similar to conventional contrastive learning Sohn (2016a); Oord et al. (2018), we employ augmented samples as

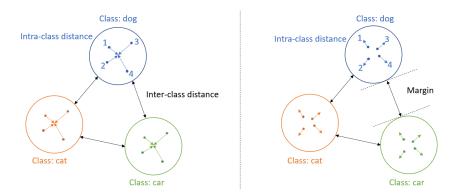


Figure 1: In most distillation methods, the teacher model is trained such that samples from different classes are separated from each other while samples from the same class are brought closer together, as illustrated in the left. However, as the intra-class distance decreases, the intra-class information within the soft labels may be lost. The proposed intra-class contrastive loss encourages an appropriate increase in intra-class distance. We also introduce a margin to ensure the stability of inter-class distances, as shown on the right.

positive samples and other samples from the same class as negative samples. This enables the teacher model to learn intra-class variability. It should be noted that, in contrast to Tian et al. (2019), they introduce contrastive learning during student training to achieve better alignment with the teacher, whereas our method employs intra-class contrastive learning during teacher training to extract the information contained within the soft labels. In practice, the intra-class contrastive loss may lead to two issues: the potential for mode collapse and low convergence speed. To address these issues, we integrate the concept of margin loss into intra-class contrastive learning. Besides, we implemented a pipeline-based caching mechanism to reduce memory usage and improve training stability under GPU memory constraints.

Furthermore, we analyzed the impact of the intra-class contrastive loss on the intra-class distances and inter-class distances. Initially, we demonstrated a quantitative relationship between intra-class contrastive loss and the distances within and between classes. Moreover, for the teacher model targeted at the proposed loss, we proved that the model's intra-class diversity and the proportion of intra-class contrastive loss satisfy specific constraints. As the proportion of intra-class contrastive loss increases, so does the intra-class diversity of the teacher model. This further substantiates the effectiveness of our proposed intra-class contrastive learning approach. Our main contributions are summarized as follows:

- We introduce an intra-class contrastive objective for training the teacher, enriching soft labels with fine-grained within-class variations. By gently dispersing samples that share the same ground-truth label, the teacher no longer collapses them to a single prototype, preventing soft labels from becoming overly deterministic and preserving inter-sample nuances that are indispensable for effective distillation.
- To counter the mode-collapse and slow-convergence issues that often plague contrastive learning, we embed a learnable margin into the intra-class contrastive loss. This margin acts as a soft barrier that keeps dispersed samples from crossing class boundaries, stabilising optimisation while retaining the discriminative structure needed for reliable knowledge transfer.
- Recognising the memory burden of large-batch contrastive learning, we propose a pipeline-style caching scheme. Features are asynchronously en-queued before back-propagation and the queue is refreshed in a sliding-window fashion, enabling the use of abundant negative samples without exceeding memory budgets and ensuring consistent contrastive signals throughout training.
- On the theoretical side, we derive an explicit relationship between the proposed loss and the feature geometry. We prove that the intra-class spread increases monotonically with the weight of the contrastive term, while the inter-class separation is guarded by the margin, thereby offering a single tunable knob that continuously controls the diversity encoded in the teacher's soft labels.

2 RELATED WORK

Knowledge Distillation. Knowledge Distillation has gained significant attention Buciluă et al. (2006); Hinton et al. (2015) in the field of deep learning. It was initially proposed for model compression Kim et al. (2018); Jin et al. (2021), and later widely adopted for knowledge transfer Vapnik et al. (2015); Zagoruyko & Komodakis (2016), or as a trick to enhance performance Guo et al. (2023); Jin et al. (2023). In traditional knowledge distillation, one teacher model teach one student model. Self-distillation Zhang et al. (2019); Lee et al. (2019) is a variant of distillation where a single model acts both as the teacher and the student. During training, a tradeoff between the ground-truth and the model's previous outputs is used as the target for subsequent training. Typically, the tradeoff parameter in self-distillation varies across epochs. In ensemble distillation You et al. (2017); Zhu et al. (2018), the outputs from multiple teachers are integrated and used to instruct the student. In mutual learning (peer learning) Zhang et al. (2018); Chen et al. (2020), multiple models learn from each other or use the aggregated knowledge as a common teacher.

Although there are many variants of distillation, a significant focus remains on the dark knowledge hidden in soft labels. Most explanations regarding dark knowledge are empirically validated. However, theoretical analyses of this concept are diverse and subject to debate. Phuong & Lampert (2019) studied the distillation mechanism with assumption that both the teacher and the student models are linear. Similarly, Allen-Zhu & Li (2020) hypothesized that data adhere to a multi-view structure, where samples from different classes may share similar features. They demonstrated the effectiveness of soft labels, which stems from their ability to enable learning of information from other classes present in the samples. Research by Mobahi et al. (2020) analyzed the teacher model in self-distillation using Green's function, and regard self-distillation as a form of regularization. On a statistical front, Menon et al. (2021) and Zhou & Song (2021) treated the generated soft labels as posterior probabilities and posited the existence of the Bayes probability. Additionally, there are other studies analyzing soft labels from the perspective of transfer risk Ji & Zhu (2020); Hsu et al. (2021).

Contrastive Learning. Contrastive Learning has emerged as a powerful mechanism for learning effective representations Logeswaran & Lee (2018); Oord et al. (2018); Tian et al. (2020a) by contrasting positive pairs against negative pairs Mikolov et al. (2013). This technique is primarily used in unsupervised or self-supervised learning environments where labeled data is scarce or expensive to obtain. Contrastive earning hinges on the idea that an encoder should map semantically similar (positive) samples closer in the embedding space, while semantically different (negative) samples should be farther apart Hadsell et al. (2006); Sohn (2016a). Sohn (2016b) introduces multiclass N-pair loss instead of traditional triplet loss. Wu et al. (2018) proposes instance discrimination to learn an embedding that can repell each pair of two training data. He et al. (2020) uses a dynamic dictionary with a momentum-updated encoder to efficiently handle large-scale data in contrastive learning so that more negative samples can be used. Gutmann & Hyvärinen (2010) assumes that the data originate from a distribution that can be parametrically clustered.

In addition to empirical work, there are some theoretical results in contrastive learning. Saunshi et al. (2019) analyzed the generalization error bounds for downstream classifiers with representations obtained from contrastive learning, based on the assumption of latent classes. Saunshi et al. (2019) argue that reducing the mutual information between views is beneficial to downstream classification accuracy. Tian et al. (2020b) proved that contrastive loss can help align features from positive pairs and features from different classes will uniformly distributed on the hypersphere. Parulekar et al. (2023) demonstrated that the representations learned by minimizing InfoNCE Loss, even with a limited number of negative samples, are consistent with the clusters inherent in the data. They further proved that when combined with a two-layer ReLU head, the learned representations can achieve zero downstream error on any binary classification task that preserves clustering. HaoChen et al. (2021) analyzed contrastive learning by constructing an augmentation graph and demonstrated that features obtained by minimizing spectral contrastive loss have provable accuracy guarantees under linear probe evaluation.

3 Method

In this section, we propose intra-class contrastive learning and define intra-class contrastive loss based on (n+1)-tuplet loss Sohn (2016b). Furthermore, we introduce margin loss to enhance model stability and accelerate convergence.

3.1 Preliminary

 Let \mathcal{X} be the sample space, $\mathcal{Y} = \{1, 2, \dots, c\}$ be the label space with c classes. Given a sample $x \in \mathcal{X}, y \in \mathcal{Y}$ is the corresponding true label. Define C(x) as the set of samples with the same class as x. In other words, $x' \in C(x)$ implies that x' and x belong to the same class. Define $\mathcal{F}: \mathcal{X} \to \mathcal{Y}$ as the hypothesis space. Each $f \in \mathcal{F}$ is a classifier. f_t and f_s are used to represent the teacher model and student model respectively. Similarly, p_t and p_s are used to represent the soft labels output by the teacher and student models. Let x^+ represent the positive sample of x. In this paper, all x^+ are augmented versions of x. Let x^- represent the negative sample of x. We employ tuplet loss as the contrastive loss function and the classical (n+1)-tuplet loss is defined as follows:

$$\mathcal{L}_{Tuplet} = \log(1 + \frac{\sum_{j=1}^{n} \exp(f(x)^{\top} f(x_j^{-}))}{\exp(f(x)^{\top} f(x_j^{+}))}). \tag{1}$$

The tuplet loss encourages reducing the distance between the positive pair (x, x^+) and increasing the distance between the negative pair (x, x^-) .

3.2 Intra-Class Contrastive Distillation

Knowledge distillation involves teaching the student model to mimic the outputs of the teacher model. The loss function for the student model is defined as:

$$\mathcal{L}_{KD} = \alpha \mathcal{L}_{CE}(y, p_s) + (1 - \alpha) \mathcal{L}_{CE}(p_t, p_s). \tag{2}$$

Here, \mathcal{L}_{CE} represents the cross-entropy loss, y denotes the ground-truth label, p_t is the soft distribution predicted by the teacher, p_s is the hard distribution predicted by the student, and α is a weighting factor. In fact, the two parts of the loss in Equation 2 can be combined into a single term: $\mathcal{L}_{KD} = \mathcal{L}_{CE}(q,p_s)$, where $q = \alpha y + (1-\alpha)p_t$. This equation provides a more intuitive demonstration of the student model's objective. From this, we can observe that distillation essentially involves making the student model approximate the soft label q. The dark knowledge in the soft label q originates from the teacher's output p_t . Existing distillation methods mostly focus on how the student aligns with the teacher, paying less attention to the learning objectives of the teacher model. In most distillation approaches, the teacher typically learns from real labels (with regularization). If the network has strong fitting capability, it may also result in the teacher's soft labels being very similar to the real labels. This can potentially lead to the reduction of intra-class information within soft labels, consequently resulting in performance degradation.

Therefore, we propose incorporating intra-class contrastive learning into teacher training to extract intra-class information. Traditional contrastive learning aims to maximize inter-class distance and minimize intra-class distance to learn discriminative and robust representations. Intra-class contrastive learning encourages the teacher model to learn embeddings where samples from the same class are dispersed appropriately, while still being distinguishable within their respective classes, thus enriching the soft label information in knowledge distillation.

In detail, we adopt contrastive loss function similar to 1. For an anchor sample x, we use its augmented view as the positive sample x^+ and other samples from the same class as negative samples x^- . The intra-class loss function is defined as:

$$\mathcal{L}_{Intra} = \log(1 + \frac{\sum_{k=1}^{m} \exp(f(x)^{\top} f(x_k^{-}))}{\exp(f(x)^{\top} f(x^{+}))}).$$
(3)

where m is the number of negative samples, $\{x_k^-\}$ is the set of negative samples. The primary distinction between the intra-class loss and the classical tuplet loss lies in the selection of negative samples. The former selects negative samples with the same class as x, while the latter mostly chooses negative samples with classes different from x. Then, the total objective function of the teacher model is defined as:

$$\mathcal{L}_{Teacher} = \mathcal{L}_{CE}(y, p_t) + \lambda \mathcal{L}_{Intra}, \tag{4}$$

Algorithm 1 SGD for Margin-Based Intra-Class Contrastive Distillation

```
217
                 Input: Training data \{(x_i, y_i)\}, learning rate \eta, margin threshold \delta, balance parameter \lambda and
218
                 maximum iteration T
219
                 Output: Trained parameters \theta
220
             1: Initialize parameters \theta of the teacher model;
221
             2: for t=1 to T do
222
             3:
                      for each batch \{x_i, y_i\} from the training data do
             4:
                           Evaluate p_{t_i} = \text{Softmax}(f_{\theta}(x_i)) for each x_i in the batch;
                          Determine \rho_{x_i} = p_{x_i}^{y_i} - \max_{j \in \mathcal{Y} \setminus \{y_i\}} p_{x_i}^j for each x_i;
             5:
224
                           Calculate \mathcal{L}_{Intra} as Eq. (5) and the total loss \mathcal{L}_{Teacher} Eq. (4);
             6:
225
                           Compute gradients \nabla_{\theta} \mathcal{L};
             7:
226
                           Update parameters \theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L};
             8:
227
                      end for
             9:
228
           10: end for
229
```

where $\lambda>0$ represents the weight of the intra-class loss. In the entire loss function, cross-entropy loss constitutes the major component. This is quite evident as the teacher model needs to first be able to distinguish between class categories before enhancing intra-class diversity. From a clustering perspective, this means that intra-class distances should be smaller than inter-class distances. By introducing intra-class loss, the representation in the teacher model becomes more enriched. This ensures that the distilled knowledge transferred to the student model is not merely a replication of the real labels but includes deeper, class-specific insights.

3.3 MARGIN-BASED INTRA-CLASS CONTRASTIVE DISTILLATION

The teacher model trained with loss 4 may encounter several issues. First, there is a partial conflict between the cross-entropy loss and the intra-class loss, as the former encourages representations of samples within the same class to be more similar. This makes the model highly sensitive to the hyperparameter λ , potentially leading to model instability. Second, for samples that the teacher cannot correctly classify, increasing the intra-class distance is not reasonable. It is not desired intra-class distances exceed inter-class distances. Third, at the beginning of training, the model should focus on learning the differences between classes, and the presence of intra-class loss can slow down convergence.

Definition 3.1. Let p^i denote the i-th degree. For each $x \in \mathcal{X}$ with the ground-truth label y, the margin is defined as $\rho_x = p_x^y - \max_{i \in \mathcal{Y} \setminus \{y\}} p_x^i$. Then, we have $-1 \le \rho_x \le 1$. A higher value of ρ_x implies a higher proportion of the ground-truth label.

We aim to select specific samples for intra-class contrastive learning through the use of margin. To achieve this, we simply choose anchor samples whose margin exceeds a specific threshold when calculating the intra-class contrastive loss.

$$\mathcal{L}_{Intra} = \begin{cases} \log(1 + \frac{\sum_{k=1}^{m} \exp(f(x)^{\top} \cdot f(x_k^{-}))}{\exp(f(x)^{\top} \cdot f(x^{+}))}) & \text{if } \rho_x > \delta \\ 0 & \text{otherwise} \end{cases}, \tag{5}$$

where $\delta>0$ is the threshold. Incorporating a margin threshold, δ , into the intra-class contrastive distillation process strategically filters the anchor samples that contribute to the loss calculation. This approach focuses on strengthening the representations of those samples which are already well classified. Besides, the teacher selectively improves the feature embeddings only for samples where the confidence margin exceeds this value, effectively ignoring those where the model's certainty is low. This results in a more stable training process as it prevents the intra-class loss from overwhelming the model with conflicting gradients from poorly classified examples. Furthermore, it ensures that the learning process is not only focused on distinguishing between classes but also on refining the model's understanding within well-understood categories, thereby enhancing the overall effectiveness of the distillation process.

In practice, we observed that due to the constraints of GPU memory capacity, the number of negative samples per class in each batch is relatively small, with even fewer meeting the threshold. This

significantly impacts the performance of the intra-class loss. To address this issue, we adopted a pipeline-based approach: samples that meet the threshold are cached in a pipeline corresponding to their class. Once the number of samples in the pipeline is sufficient, we compute the intra-class loss and then clear the pipeline. This method not only substantially reduces memory consumption but also enhances the stability of intra-class contrastive learning.

In this section, we introduced the intra-Class Contrastive loss and discussed how it assists in generating soft labels by the teacher model. Addressing potential issues with the proposed intra-Class Contrastive Distillation, we have incorporated the concept of margin to improve the intra-Class Contrastive loss. The complete algorithm for the teacher model can be found in Algorithm 1.

4 THEORETICAL ANALYSIS

4.1 FORMULATION

In this section, we analyze the representations learned by the teacher model from the perspective of clustering. Our focus here is not on the model's classification performance, but rather on the distances within and between classes. Thus, in this section, the teacher model works as a feature embedding function $\varphi: \mathcal{X} \to \mathbb{R}^d (d > c)$, which transforms the data point from the m-dimensional sample space to the d-dimensional embedding space. We assume that φ is normalized, such that $\|\varphi(x)\|_2 = 1$ for any $x \in \mathcal{X}$. Denote \mathcal{H} as the hypothesis space of all embedding functions. For any sample x, the positive sample x^+ is an augmented version of x, while the negative samples x^- are obtained through sampling.

In the previous sections, we proposed a method based on a tradeoff between cross-entropy loss and intra-class contrastive loss. In this section, we analyze the impact of the intra-class contrastive loss on the model. Since the teacher model performs supervised learning (with visible ground-truth labels), we employ cross-entropy loss to effectively learn inter-class differences. However, for the convenience of the following theoretical analysis, we substitute the loss function with a tradeoff between conventional inter-class contrastive loss and intra-class contrastive loss, i.e.,

$$\mathcal{L}(\varphi) = \underbrace{\log(1 + \frac{\sum_{j=1}^{n} \exp\left(\varphi(x)^{T} \varphi(x_{j}^{-})\right)}{\exp\left(\varphi(x)^{T} \varphi(x^{+})\right)}}_{\mathcal{L}_{Inter}}) + \lambda \cdot \underbrace{\left(\log(1 + \frac{\sum_{k=1}^{m} \exp\left(\varphi(x)^{T} \varphi(x_{k}^{-})\right)}{\exp\left(\varphi(x)^{T} \varphi(x^{+})\right)}\right)}_{\mathcal{L}_{Intra}}\right)}_{\mathcal{L}_{Intra}}.$$
(6)

Here, x^+ is an augmented version of the sample x, serving as the positive sample, while x_j^- and x_k^- are negative samples. Notably, $x_j^- \notin C(x)$ in $\mathcal{L}_{\text{Inter}}$, whereas $x_k^- \in C(x)$ in $\mathcal{L}_{\text{Intra}}$. The inter-class loss $\mathcal{L}_{\text{Inter}}$ enables the model to learn the differences between different classes, ensuring that samples from different categories are well-separated in the feature space. Conversely, the intra-class loss $\mathcal{L}_{\text{Intra}}$ focuses on learning the differences within the same class, ensuring that samples of the same category are appropriately mutually distant in the feature space.

4.2 Inter-class distances and intra-class distances

Before investigating our proposed method, we first consider a teacher model trained directly using cross-entropy loss.

Proposition 4.1. Consider $\lambda = 0$. If the model has perfect fitting capabilities, then the teacher model $f_T = \arg\min \mathcal{L}_{Teacher}$ in Eq.4 will induce soft labels that are identical to the ground-truth labels.

This conclusion is evident. It shows that without some guidance for the teacher model, the generated soft labels may lose a significant amount of intra-class information. Now, we focus on the specifics of how our approach manages distance metrics within the embedding space. Essentially, the proposed inter-class contrastive loss and intra-class contrastive loss are designed to control inter-class distances and intra-class distances respectively. First, we define them as follows. Given the embedding φ , the inter-class distance is defined as

$$d_{\text{Inter}} = \mathbb{E}_{x^- \notin C(x)} \left[\exp(\varphi(x)^T \varphi(x^-)) \right],$$

and the intra-class distance is defined as

$$d_{\text{Intra}} = \mathbb{E}_{x^- \in C(x)} \left[\exp(\varphi(x)^T \varphi(x^-)) \right].$$

Table 1: **Results on CIFAR-100 and Tiny ImageNet, Homogeneous Architecture.** Top-1 accuracy is adopted as the evaluation criterion. All experiments are repeated 5 times, and the table presents the final mean of the results. The best results are presented in bold.

		CIFAR-100		Tiny ImageNet			
Teacher	ResNet50	WRN-40-2	VGG13	ResNet50	WRN-40-2	VGG13	
Teacher	78.31	76.89	74.40	70.57	71.34	67.23	
Student	ResNet34	WRN-16-2	VGG8	ResNet34	WRN-16-2	VGG8	
Student	78.19	76.40	73.80	67.14	67.60	63.60	
KD	78.38	76.60	73.99	67.33	68.83	64.00	
FitNet	75.20	75.03	72.33	66.96	67.74	57.67	
RKD	78.65	78.30	74.73	68.47	69.23	63.25	
CRD	78.57	78.69	74.48	69.45	70.27	64.59	
OFD	76.64	76.88	72.64	66.35	66.12	58.79	
ReviewKD	77.55	77.94	73.24	68.02	69.45	63.57	
VID	77.43	77.63	73.51	67.89	68.27	62.67	
MLLD	79.09	79.50	75.07	69.48	70.78	64.82	
Ours	79.10	79.09	74.96	70.22	71.09	65.14	
Ours+RKD	79.39	79.53	75.70	70.96	72.08	66.32	

Theorem 4.2. As the two numbers of negative samples in 6 $n, m \to \infty$ and n/m = K, for a certain φ , we have $\frac{d_{lntra}}{d_{lnter}} = K \frac{\exp(\mathcal{L}_{lntra})}{\exp(\mathcal{L}_{lnter})}$.

In fact, Theorem 4.2 connects the loss defined in 6 with inter-class and intra-class distances, illustrating that the proposed loss can effectively control both intra-class and inter-class distances. Furthermore, although the theorem requires that n and m approach infinity at the same rate, in practice, we can use the ratio of losses with finite samples as an approximation. Given the framework established by the previous sections and particularly by Theorem 4.2, we can explore the practical implications of these findings and further justify our method's approach. The two theorems essentially state that the relationship between the intra-class and inter-class distances can be quantitatively managed through the loss ratio.

Next, we will consider the impact of the intra-class contrastive loss on the model when optimizing the objective function

$$\min_{\varphi \in \mathcal{H}} \left\{ \mathcal{L}(\varphi) = \mathcal{L}_{\text{Inter}}(\varphi) + \lambda \mathcal{L}_{\text{Intra}}(\varphi) \right\}. \tag{7}$$

Theorem 4.3. Assume that $\varphi^* \in \arg \min_{\varphi \in \mathcal{H}} \mathcal{L}(\varphi)$, and then we have

$$\frac{1}{C_0 \cdot \lambda + C_1} \le \frac{\mathcal{L}_{Intra}}{\mathcal{L}_{Inter}} \le C_2 \cdot \frac{1}{\lambda} + C_3, \tag{8}$$

where C_0, C_1, C_2, C_3 are positive constants about m and n.

This theorem provides bounds for the ratio of intra-class and inter-class losses when optimizing the objective function $\mathcal{L}(\varphi)$. It shows how the balance parameter λ affects the ratio between $\mathcal{L}_{\text{Inter}}$ and $\mathcal{L}_{\text{Intra}}$. The lower and upper bounds of this ratio are inversely related to λ . Thus, by adjusting λ , one can effectively control the balance between intra-class and inter-class optimization. A larger λ leads to tighter intra-class samples, while a smaller λ promotes better separation between different classes.

These theoretical results provide a solid justification for our approach by demonstrating how the contrastive loss, when modulated with a carefully chosen balance parameter λ , effectively manages the trade-off between intra-class compactness and inter-class separation. Moreover, the bounds in Theorem 4.3 offer assurance for both the control of intra-class diversity and the stability of the training process.

Table 2: **Results on CIFAR-100 and Tiny ImageNet, Heterogeneous Architecture.** Top-1 accuracy is adopted as the evaluation criterion. All experiments are repeated 5 times, and the table presents the final mean of the results. The best results are presented in bold.

	CIFAR-100			Tiny ImageNet			
Teacher	ResNet50	VGG13	WRN-40-2	ResNet50	VGG13	WRN-40-2	
Teacher	78.31	74.40	76.89	70.57	67.23	71.34	
Student	MobileNet-V2	MobileNet-V2	ShuffleNet-V2	MobileNet-V2	MobileNet-V2	ShuffleNet-V2	
Student	65.18	65.18	69.23	50.76	50.86	53.37	
KD	66.22	65.82	70.16	50.39	52.87	55.77	
FitNet	67.52	63.51	65.06	54.37	51.81	46.53	
RKD	65.41	66.51	72.11	55.58	56.27	59.58	
CRD	65.79	66.46	71.08	57.21	58.31	61.24	
ReviewKD	67.59	65.67	72.42	56.44	55.97	60.35	
VID	67.88	66.18	72.71	54.98	56.08	59.44	
MLLD	68.99	65.58	65.89	53.72	54.13	58.76	
Ours	68.72	66.44	72.00	57.31	59.12	61.78	
Ours+RKD	69.03	66.79	72.86	58.12	60.09	63.14	

The choice of λ . From a practical standpoint, Theorem 4.3 illustrates the tradeoff between the two losses with respect to λ . In conjunction with Theorem 4.2, it allows us to adjust λ to achieve an optimal balance between intra-class and inter-class distances. By choosing an appropriate value of λ , it is ensured that the model does not bias too much towards distinguishing only between classes and neglecting the variance within the classes, or vice versa. This balance is particularly beneficial in applications where subtle intra-class variations are as significant as the differences between classes.

In conclusion, the theoretical insights provided by the analysis not only bolster the validity of using a margin-based intra-class contrastive distillation approach but also highlight the importance of carefully considering the balance of loss components to achieve the best learning outcomes. As we move forward, these principles can guide the development of more sophisticated models that are tuned to the nuances of specific tasks and data characteristics.

5 Experiments

In this section, we evaluate the performance of the proposed Margin-Based Intra-Class Contrastive Distillation algorithm on image classification datasets. We assessed the effectiveness of distillation across different model architectures.

5.1 Datasets and Settings

 Setting. We compared two common settings in knowledge distillation: (1) Homogeneous architecture, where the teacher and student models share the same architecture, and (2) Heterogeneous architecture, where the teacher and student models have different architectures. Additionally, for each setting, we conducted experiments for various neural network architectures. The experimental results presented are the averages from 5 repeated trials. Owing to the limitations of page width, we have presented only the mean values without the standard deviation.

Baselines. To validate the effectiveness of our proposed method, we conducted experiments comparing it with the vanilla KD Hinton et al. (2015) and seven benchmark methods: FitNet Romero et al. (2014), RKD Park et al. (2019), CRD Tian et al. (2019), OFD Heo et al. (2019), ReviewKD Chen et al. (2021b), VID Ahn et al. (2019) and MLLD Jin et al. (2023). These methods provide a comprehensive backdrop against which the performance of our approach can be measured.

Unlike many traditional distillation methods that primarily focus on aligning the student model with the teacher model, our proposed algorithm specifically targets the teacher model to enrich the information contained in the soft labels. This distinction is crucial because soft labels are a

Table 3: **Results on ImageNet.** Top-1 and Top-5 accuracy is adopted as the evaluation metric. The best results are presented in bold.

Student (Teacher)	Metric	Teacher	Student	KD	RKD	CRD	ReviewKD	MLLD	Ours
ResNet-18 (ResNet-34)	Top-1	68.53	64.95	66.83	66.52	65.69	67.29	65.52	67.81
	Top-5	86.44	83.68	84.39	85.12	85.15	84.60	84.89	89.34
MobileNet-V2 (ResNet-50)	Top-1	72.77	60.77	62.93	67.27	65.61	67.22	68.90	69.57
	Top-5	88.18	84.14	83.44	84.69	86.84	87.34	86.51	88.67

 critical element in guiding the student model during distillation. By improving the quality of the teacher model's soft labels, our approach allows the student model to learn more nuanced information. Consequently, our method can be integrated with many existing distillation techniques to further enhance performance. For all experiments conducted, we compare our Margin-Based Intra-Class Contrastive Distillation (ours) with other standard distillation algorithms. Additionally, we combined our algorithm with the classical RKD Park et al. (2019) (ours+RKD) to evaluate whether this combination results in even better performance.

Experimental Results and Analysis. For CIFAR-100 and Tiny ImageNet, we conducted six experimental groups under both homogeneous and heterogeneous settings. Due to space limitations, we present three representative results in Table 1 and Table 2, respectively. The remaining results are provided in the appendix (see Table 4 and 5). The results for ImageNet are reported in Table 3. Even without integrating with RKD, our method, which relies solely on the enhanced soft labels for conventional distillation, nearly surpasses other methods. This further validates the effectiveness of extracting information from the soft labels. Moreover, it illustrates that intra-class contrastive loss facilitates the learning of improved soft labels, thereby strengthening the teacher model's teaching capability and enhancing distillation effectiveness. We plot the T-SNE on CIFAR-100 (see Figure 2 in

Ablation Study and Hyperparameter Sensitivity Analysis. We conduct comprehensive ablation studies to evaluate the effectiveness of the proposed margin loss in both homogeneous and heterogeneous architectures, as shown in Tables 8 and 9 in the appendix. In the homogeneous setting (Table 8), where the teacher and student share similar architectures, adding margin loss consistently improves student model performance across both CIFAR-100 and Tiny ImageNet datasets. The trend remains consistent in the heterogeneous setting (Table 9). To further assess the stability of our method, we investigate the sensitivity to the hyperparameter λ , which controls the strength of the margin loss. The result is shown in Figure 3 in the appendix.

the appendix) and it verifies intra-class contrastive loss can increase intra-class diversity. Additionally,

our approach is compatible with many existing distillation algorithms, further improving performance.

Training time. We compared the training time of the teacher model with intra-class contrastive loss and the standard teacher model. The results are presented in Table 7 in the appendix. We conducted the comparison by evaluating the training duration for each data batch. The results show that after implementing the pipeline-based caching mechanism, the additional computational overhead amounts to approximately 10%-15%.

6 Conclusion

In this work, we introduced the Margin-Based Intra-Class Contrastive Distillation approach, which integrates intra-class contrastive learning with traditional knowledge distillation. This method not only enriches the soft labels with nuanced intra-class variations but also maintains a balance between inter-class and intra-class distinctions, which is vital for the robust generalization of the student model. Our method significantly enhances the performance of the student model by leveraging enriched soft labels, demonstrating superior results across standard image classification datasets. Both the integration of margin loss and the design of pipeline ensure the stability and efficiency of the learning process, addressing potential issues related to convergence and model training dynamics. The theoretical results also confirm the feasibility of our method and highlight the role of the parameter λ in balancing the tradeoff. Overall, our approach provides a compelling framework for effectively utilizing the soft labels in knowledge distillation, paving the way for future innovations in model compression and efficient learning strategies. A limitation of our method is that the weight of the intra-class contrastive loss needs to be adjusted depending on the dataset.

REFERENCES

- Sungsoo Ahn, Shell Xu Hu, Andreas Damianou, Neil D Lawrence, and Zhenwen Dai. Variational information distillation for knowledge transfer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9163–9171, 2019.
- Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *CoRR*, abs/2012.09816, 2020.
- 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.
- Defang Chen, Jian-Ping Mei, Can Wang, Yan Feng, and Chun Chen. Online knowledge distillation with diverse peers. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 3430–3437, 2020.
- Liqun Chen, Dong Wang, Zhe Gan, Jingjing Liu, Ricardo Henao, and Lawrence Carin. Wasserstein contrastive representation distillation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16296–16305, 2021a.
- Pengguang Chen, Shu Liu, Hengshuang Zhao, and Jiaya Jia. Distilling knowledge via knowledge review. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5008–5017, 2021b.
- Wenliang Dai, Lu Hou, Lifeng Shang, Xin Jiang, Qun Liu, and Pascale Fung. Enabling multimodal generation on clip via vision-language knowledge distillation. *arXiv preprint arXiv:2203.06386*, 2022.
- 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.
- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models. In *The Twelfth International Conference on Learning Representations*, 2023.
- Ziyao Guo, Haonan Yan, Hui Li, and Xiaodong Lin. Class attention transfer based knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11868–11877, 2023.
- Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In Yee Whye Teh and Mike Titterington (eds.), *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pp. 297–304, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR.
- Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In 2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06), volume 2, pp. 1735–1742. IEEE, 2006.
- Jeff Z HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for self-supervised deep learning with spectral contrastive loss. *Advances in Neural Information Processing Systems*, 34:5000–5011, 2021.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9729–9738, 2020.
- Byeongho Heo, Jeesoo Kim, Sangdoo Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A comprehensive overhaul of feature distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1921–1930, 2019.
 - Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531, 2015.

- Daniel Hsu, Ziwei Ji, Matus Telgarsky, and Lan Wang. Generalization bounds via distillation. *arXiv* preprint arXiv:2104.05641, 2021.
- Guangda Ji and Zhanxing Zhu. Knowledge distillation in wide neural networks: Risk bound, data efficiency and imperfect teacher. *Advances in Neural Information Processing Systems*, 33: 20823–20833, 2020.
 - Zhong Ji, Jingwei Ni, Xiyao Liu, and Yanwei Pang. Teachers cooperation: team-knowledge distillation for multiple cross-domain few-shot learning. *Frontiers of Computer Science*, 17(2):172312, 2023.
 - Qing Jin, Jian Ren, Oliver J Woodford, Jiazhuo Wang, Geng Yuan, Yanzhi Wang, and Sergey Tulyakov. Teachers do more than teach: Compressing image-to-image models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13600–13611, 2021.
 - Ying Jin, Jiaqi Wang, and Dahua Lin. Multi-level logit distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 24276–24285, 2023.
 - Jangho Kim, SeongUk Park, and Nojun Kwak. Paraphrasing complex network: Network compression via factor transfer. *Advances in neural information processing systems*, 31, 2018.
 - Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
 - Yann Le and Xuan Yang. Tiny imagenet visual recognition challenge. CS 231N, 7(7):3, 2015.
 - Hankook Lee, Sung Ju Hwang, and Jinwoo Shin. Rethinking data augmentation: Self-supervision and self-distillation. *CoRR*, abs/1910.05872, 2019.
 - Lajanugen Logeswaran and Honglak Lee. An efficient framework for learning sentence representations. *arXiv preprint arXiv:1803.02893*, 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.
 - Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26, 2013.
 - Hossein Mobahi, Mehrdad Farajtabar, and Peter Bartlett. Self-distillation amplifies regularization in hilbert space. *Advances in Neural Information Processing Systems*, 33:3351–3361, 2020.
 - Alexander Neitz, Giambattista Parascandolo, and Bernhard Schölkopf. A teacher-student framework to distill future trajectories. In *International Conference on Learning Representations*, 2020.
 - Mehdi Noroozi, Ananth Vinjimoor, Paolo Favaro, and Hamed Pirsiavash. Boosting self-supervised learning via knowledge transfer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9359–9367, 2018.
 - Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
 - Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3967–3976, 2019.
 - Advait Parulekar, Liam Collins, Karthikeyan Shanmugam, Aryan Mokhtari, and Sanjay Shakkottai. Infonce loss provably learns cluster-preserving representations. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 1914–1961. PMLR, 2023.
 - Mary Phuong and Christoph Lampert. Towards understanding knowledge distillation. In *International Conference on Machine Learning*, pp. 5142–5151. PMLR, 2019.
 - Ahmad Rashid, Vasileios Lioutas, Abbas Ghaddar, and Mehdi Rezagholizadeh. Towards zero-shot knowledge distillation for natural language processing. *arXiv preprint arXiv:2012.15495*, 2020.

- 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.
 - Nikunj Saunshi, Orestis Plevrakis, Sanjeev Arora, Mikhail Khodak, and Hrishikesh Khandeparkar. A theoretical analysis of contrastive unsupervised representation learning. In *International Conference on Machine Learning*, pp. 5628–5637. PMLR, 2019.
 - Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. *Advances in neural information processing systems*, 29, 2016a.
 - Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. *Advances in neural information processing systems*, 29, 2016b.
 - Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. In *International Conference on Learning Representations*, 2019.
 - Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*, pp. 776–794. Springer, 2020a.
 - Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What makes for good views for contrastive learning? *Advances in neural information processing systems*, 33:6827–6839, 2020b.
 - Vladimir Vapnik, Rauf Izmailov, et al. Learning using privileged information: similarity control and knowledge transfer. *J. Mach. Learn. Res.*, 16(1):2023–2049, 2015.
 - Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3733–3742, 2018.
 - Ning Xu, Yihao Hu, Congyu Qiao, and Xin Geng. Aligned objective for soft-pseudo-label generation in supervised learning. In *Forty-first International Conference on Machine Learning*.
 - Ziqing Yang, Yiming Cui, Zhipeng Chen, Wanxiang Che, Ting Liu, Shijin Wang, and Guoping Hu. Textbrewer: An open-source knowledge distillation toolkit for natural language processing. *arXiv* preprint arXiv:2002.12620, 2020.
 - Shan You, Chang Xu, Chao Xu, and Dacheng Tao. Learning from multiple teacher networks. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1285–1294, 2017.
 - Fei Yuan, Linjun Shou, Jian Pei, Wutao Lin, Ming Gong, Yan Fu, and Daxin Jiang. Reinforced multi-teacher selection for knowledge distillation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 14284–14291, 2021.
 - Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. abs/1612.03928, 2016.
 - Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3713–3722, 2019.
 - Ying Zhang, Tao Xiang, Timothy M Hospedales, and Huchuan Lu. Deep mutual learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4320–4328, 2018.
 - Helong Zhou and Liangchen Song. Rethinking soft labels for knowledge distillation: A bias-variance tradeoff perspective. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2021.
 - Xiatian Zhu, Shaogang Gong, et al. Knowledge distillation by on-the-fly native ensemble. *Advances in neural information processing systems*, 31, 2018.

A PROOF OF THEOREM 4.2

By the definition of inter-class distance and intra-class distance, given a certain φ , we have

$$\frac{d_{\text{Intra}}}{d_{\text{Inter}}} = \frac{\mathbb{E}_{x^{-} \notin C(x)} \left[e^{\varphi(x)^{T} \varphi(x^{-})} \right]}{\mathbb{E}_{x^{-} \in C(x)} \left[e^{\varphi(x)^{T} \varphi(x^{-})} \right]}$$

$$= \lim_{\substack{m, n \to \infty \\ \frac{n}{m} = K}} \frac{\frac{1}{m} \sum_{i=1}^{m} e^{\varphi(x)^{T} \varphi(x^{-})}}{\frac{1}{n} \sum_{j=1}^{n} e^{\varphi(x)^{T} \varphi(x^{-})}}$$

$$= K \lim_{\substack{m, n \to \infty \\ \frac{n}{m} = K}} \frac{\sum_{i=1}^{m} e^{\varphi(x)^{T} \varphi(x^{-})}}{\sum_{j=1}^{n} e^{\varphi(x)^{T} \varphi(x^{-})}}.$$
(9)

On the other side, according to the definitions of \mathcal{L}_{Inter} and \mathcal{L}_{Intra} , there exists a relationship

$$\frac{\sum_{i=1}^{m} e^{\varphi(x)^{T} \varphi(x^{-})}}{\sum_{i=1}^{n} e^{\varphi(x)^{T} \varphi(x^{-})}} = \frac{e^{\mathcal{L}_{\text{Intra}}} - 1}{e^{\mathcal{L}_{\text{Inter}}} - 1}.$$
(10)

Then, when $n, m \to \infty$ and n/m = K, the constant 1 can be ignored, and then we obtain

$$\frac{d_{\text{Intra}}}{d_{\text{Inter}}} = K \frac{e^{\mathcal{L}_{\text{Intra}}}}{e^{\mathcal{L}_{\text{Inter}}}}.$$
(11)

B PROOF OF THEOREM 4.3

Recall that the total loss $\mathcal{L}(\varphi) = \mathcal{L}_{Inter}(\varphi) + \lambda \mathcal{L}_{Intra}(\varphi)$. Define

$$\mathcal{L}_{Inter}^{0}(\varphi) = \mathcal{L}_{Inter}(\varphi) - \min_{\varphi \in \mathcal{H}} \mathcal{L}_{Inter}(\varphi)$$
 (12)

$$\mathcal{L}_{Intra}^{0}(\varphi) = \mathcal{L}_{Intra}(\varphi) - \min_{\varphi \in \mathcal{H}} \mathcal{L}_{Intra}(\varphi)$$
(13)

$$\mathcal{L}^{0}(\varphi) = \mathcal{L}_{Inter}^{0} + \lambda \mathcal{L}_{Intra}^{0}.$$
 (14)

In fact, since \mathcal{L}_{Inter} and \mathcal{L}_{Intra} are both non-negative, we can approximate the lower bound with an arbitrarily small ε gap. For simplicity, we assume that the lower bound can be achieved, and let $\varphi_1 = \arg\min_{\varphi \in \mathcal{H}} \mathcal{L}_{Inter}$ and $\varphi_2 = \arg\min_{\varphi \in \mathcal{H}} \mathcal{L}_{Intra}$. Thus, it is clear that $\mathcal{L}^0_{Inter}(\varphi_1) = 0$ and $\mathcal{L}^0_{Intra}(\varphi_2) = 0$ hold.

Next, we provide a precise estimate for the lower bound of $\min_{\varphi \in \mathcal{H}} \mathcal{L}_{Inter}(\varphi)$ and $\min_{\varphi \in \mathcal{H}} \mathcal{L}_{Intra}(\varphi)$.

$$\mathcal{L}_{Inter}(\varphi) = \log \left(1 + \frac{\sum_{j=1}^{n} e^{\varphi(x)^{T}} \varphi(x_{j}^{-})}{e^{\varphi(x)^{T}} \varphi(x^{+})} \right)$$

$$\geq \log \left(1 + \frac{\sum_{j=1}^{n} e^{-1}}{e^{1}} \right)$$

$$= \log \left(1 + ne^{-2} \right).$$
(15)

The equality is satisfied when $\varphi(x)^T \varphi(x^+) = 1$ and $\varphi(x)^T \varphi(x_j^-) = -1$, for $j = 1, 2, \dots, n$. On the other side,

$$\mathcal{L}_{Intra}(\varphi) = \log \left(1 + \frac{\sum_{i=1}^{m} e^{\varphi(x)^{T}} \varphi(x_{i}^{-})}{e^{\varphi(x)^{T}} \varphi(x^{+})} \right)$$

$$\geq \log \left(1 + \frac{\sum_{i=1}^{m} e^{-1}}{e^{1}} \right)$$

$$= \log \left(1 + me^{-2} \right).$$
(16)

The equality is satisfied when $\varphi(x)^T \varphi(x^+) = 1$ and $\varphi(x)^T \varphi(x_i^-) = -1$, for $i = 1, 2, \dots, m$.

Assume that $\varphi^* \in \arg\min_{\varphi \in \mathcal{H}} \mathcal{L}(\varphi)$. Thus, it is evident that $\varphi^* \in \arg\min_{\varphi \in \mathcal{H}} \mathcal{L}^0(\varphi)$. Since that $\mathcal{L}^0_{\text{Inter}}(\varphi_1) = 0$, by the optimality of φ^* , we have that

$$\mathcal{L}_{Inter}^{0}(\varphi^{*}) + \lambda \mathcal{L}_{Intra}^{0}(\varphi^{*}) \leq \mathcal{L}_{Inter}^{0}(\varphi_{1}) + \lambda \mathcal{L}_{Intra}^{0}(\varphi_{1})$$

$$= \lambda \mathcal{L}_{Intra}^{0}(\varphi_{1}).$$
(17)

Upon expanding the equation, we obtain

$$\mathcal{L}_{Inter}(\varphi^*) - \min_{\varphi \in \mathcal{H}} \mathcal{L}_{Inter}(\varphi) + \lambda \mathcal{L}_{Intra}(\varphi^*) - \lambda \min_{\varphi \in \mathcal{H}} \mathcal{L}_{Intra}(\varphi)$$

$$\leq \lambda \mathcal{L}_{Intra}(\varphi_1) - \lambda \min_{\varphi \in \mathcal{H}} \mathcal{L}_{Intra}(\varphi).$$
(18)

Rearrange this equation, and thus

$$\frac{\mathcal{L}_{Intra}(\varphi^{*})}{\mathcal{L}_{Intra}(\varphi^{*})} \leq \frac{\lambda \mathcal{L}_{Intra}(\varphi_{1}) - \lambda \mathcal{L}_{Intra}(\varphi^{*}) + \min_{\varphi \in \mathcal{H}} \mathcal{L}_{Inter}(\varphi)}{\mathcal{L}_{Intra}(\varphi^{*})} \\
= \lambda \frac{\mathcal{L}_{Intra}(\varphi_{1})}{\mathcal{L}_{Intra}(\varphi^{*})} + \frac{\min_{\varphi \in \mathcal{H}} \mathcal{L}_{Inter}(\varphi)}{\mathcal{L}_{Intra}(\varphi^{*})} - \lambda \\
= \lambda \frac{\log\left(1 + \frac{\sum_{i=1}^{m} e^{\varphi_{1}(x)^{T}} \varphi_{1}(x_{i}^{-})}{e^{\varphi_{1}(x)^{T}} \varphi_{1}(x_{i}^{-})}\right)}{\log\left(1 + \frac{\sum_{i=1}^{m} e^{\varphi^{*}(x)^{T}} \varphi^{*}(x_{i}^{-})}{e^{\varphi^{*}(x)^{T}} \varphi^{*}(x_{i}^{+})}\right)} + \frac{\log\left(1 + ne^{-2}\right)}{\log\left(1 + \frac{\sum_{i=1}^{m} e^{\varphi^{*}(x)^{T}} \varphi^{*}(x_{i}^{-})}{e^{\varphi^{*}(x)^{T}} \varphi^{*}(x_{i}^{+})}\right)} - \lambda \\
\leq \lambda \frac{\log\left(1 + \frac{me}{e^{-1}}\right)}{\log\left(1 + \frac{me^{-1}}{e}\right)} + \frac{\log\left(1 + ne^{-2}\right)}{\log\left(1 + \frac{me^{-1}}{e}\right)} - \lambda \\
= \left(\frac{\log\left(1 + me^{2}\right)}{\log\left(1 + me^{2}\right)} - 1\right)\lambda + \frac{\log\left(1 + ne^{-2}\right)}{\log\left(1 + me^{-2}\right)} \\
= C_{0} \cdot \lambda + C_{1}, \tag{19}$$

where $C_0, C_1 > 0$.

Since both sides are non-negative, taking the reciprocal yields

$$\frac{\mathcal{L}_{Intra}(\varphi^*)}{\mathcal{L}_{Inter}(\varphi^*)} \ge \frac{1}{C_0 \cdot \lambda + C_1}.$$
 (20)

Similarly, for φ_2 satisfying $\mathcal{L}^0_{Intra}(\varphi_2) = 0$, by the optimality of φ^* , we have that

$$\mathcal{L}_{Inter}^{0}(\varphi^{*}) + \lambda \mathcal{L}_{Intra}^{0}(\varphi^{*}) \leq \mathcal{L}_{Inter}^{0}(\varphi_{2}) + \lambda \mathcal{L}_{Intra}^{0}(\varphi_{2})$$

$$= \mathcal{L}_{Inter}^{0}(\varphi_{2}). \tag{21}$$

Upon expanding the equation, we obtain

$$\mathcal{L}_{Inter}(\varphi^*) - \min_{\varphi \in \mathcal{H}} \mathcal{L}_{Inter}(\varphi) + \lambda \mathcal{L}_{Intra}(\varphi^*) - \lambda \min_{\varphi \in \mathcal{H}} \mathcal{L}_{Intra}(\varphi)$$

$$\leq \mathcal{L}_{Inter}(\varphi_2) - \min_{\varphi \in \mathcal{H}} \mathcal{L}_{Inter}(\varphi).$$
(22)

Table 4: Additional Results on CIFAR-100 and Tiny ImageNet, Homogeneous Architecture. Top-1 accuracy is adopted as the evaluation criterion. The best results are presented in bold.

		CIFAR-100		Tiny ImageNet			
Teacher	ResNet101	ResNet32×4	WRN-40-2	ResNet101	ResNet32×4	WRN-40-2	
reaction	79.17	79.86	76.89	72.47	72.91	71.34	
Student	ResNet34	ResNet8×4	WRN-40-1	ResNet34	ResNet8×4	WRN-40-1	
Student	78.19	72.92	72.23	67.14	64.60	65.12	
KD	78.31	73.55	74.02	67.99	65.11	66.64	
FitNet	77.63	74.22	73.74	67.01	66.16	67.61	
RKD	79.94	74.02	74.20	69.03	66.02	68.03	
CRD	80.21	73.26	73.07	69.29	65.03	66.88	
OFD	78.36	74.96	73.49	69.10	66.92	66.58	
ReviewKD	79.47	72.12	74.13	69.73	67.38	67.37	
VID	79.21	72.32	74.26	69.49	65.08	67.70	
MLLD	80.05	73.47	72.29	70.43	67.19	65.91	
Ours	80.70	75.44	74.54	71.81	68.30	68.48	
Ours+RKD	81.15	76.22	75.31	72.01	69.56	69.25	

Table 5: Additional Results on CIFAR-100 and Tiny ImageNet, Heterogeneous Architecture. Top-1 accuracy is adopted as the evaluation criterion. The best results are presented in bold.

		CIFAR-100		Tiny ImageNet			
Teacher	ResNet32×4	WRN-40-2	ResNet32×4	ResNet32×4	WRN-40-2	ResNet32×4	
reaction	79.86	76.89	79.86	72.91	71.34	72.91	
Student	MobileNet-V2	MobileNet-V2	ShuffleNet-V2	MobileNet-V2	MobileNet-V2	ShuffleNet-V2	
Student	65.18	65.18	69.23	50.76	50.76	53.37	
KD	66.27	65.90	70.64	51.55	52.68	55.92	
FitNet	65.72	65.00	69.18	53.31	54.13	57.20	
RKD	67.87	68.45	72.41	54.17	53.70	58.42	
CRD	67.52	67.97	72.34	53.40	53.15	59.70	
ReviewKD	68.01	66.94	72.47	53.91	53.36	59.48	
VID	67.91	67.89	73.11	54.00	51.49	58.91	
MLLD	68.43	67.05	69.07	55.42	52.74	59.20	
Ours	69.21	68.15	73.10	56.54	55.60	60.96	
Ours+RKD	69.32	68.95	73.88	57.88	56.19	62.00	

Rearrange this equation, and thus

$$\frac{\mathcal{L}_{Intra}(\varphi^{*})}{\mathcal{L}_{Inter}(\varphi^{*})} \leq \frac{\mathcal{L}_{Inter}(\varphi_{2}) + \lambda \min_{\varphi \in \mathcal{H}} \mathcal{L}_{Intra}(\varphi)}{\lambda \mathcal{L}_{Inter}(\varphi^{*})} - \frac{1}{\lambda}$$

$$= \frac{\log \left(1 + \frac{\sum_{j=1}^{n} e^{\varphi_{2}(x)^{T}} \varphi_{2}(x_{j}^{-})}{e^{\varphi_{2}(x)^{T}} \varphi_{2}(x^{+})}\right) + \lambda \log(1 + me^{-2})}{\lambda \log \left(1 + \frac{\sum_{j=1}^{n} e^{\varphi^{*}(x)^{T}} \varphi^{*}(x_{j}^{-})}{e^{\varphi^{*}(x)^{T}} \varphi^{*}(x^{+})}\right)} - \frac{1}{\lambda}$$

$$\leq \frac{\log \left(1 + \frac{ne}{e^{-1}}\right)}{\lambda \log \left(1 + \frac{ne^{-1}}{e}\right)} - \frac{1}{\lambda} + \frac{\log(1 + me^{-2})}{\log \left(1 + \frac{ne^{-1}}{e}\right)}$$

$$= \left(\frac{\log \left(1 + ne^{2}\right)}{\log \left(1 + ne^{2}\right)} - \frac{15}{\lambda}\right) \frac{1}{\lambda} + \frac{\log(1 + me^{-2})}{\log \left(1 + ne^{-2}\right)}$$

Table 6: **Results on CIFAR-10, Homogeneous and Heterogeneous Architectures.** Top-1 accuracy is adopted as the evaluation criterion. All experiments are repeated 5 times, and the table presents the final mean of the results. The best results are presented in bold.

		CIFAR	R-10 Homogen	eous	CIFAR-10 Heterogeneous			
	Teacher	ResNet50	WRN-40-2	VGG13	ResNet50	WRN-40-2	VGG13	
Method	Teacher	94.08	94.85	91.54	94.08	94.85	91.54	
Method	Student	ResNet34	WRN-16-2	VGG8	MobileNet-V2	ShuffleNet-V2	MobileNet-V2	
	Student	93.16	93.22	91.04	83.46	87.09	83.46	
K	D D	93.92	93.82	91.96	84.72	88.64	78.83	
Fit	Net	93.53	93.74	91.38	86.39	85.48	82.70	
RI	KD	94.13	94.05	92.22	86.26	89.41	85.21	
CI	RD	90.97	90.70	89.31	85.43	88.12	83.15	
Ol	FD	94.08	94.09	92.01	85.76	89.02	81.83	
Revie	ewKD	94.04	94.02	92.18	86.09	90.86	83.97	
V	ID	94.02	93.99	91.87	86.49	90.31	86.23	
ML	LLD	92.32	92.18	90.28	85.78	87.33	84.12	
O	urs	94.17	94.23	92.33	86.51	90.96	85.98	
Ours-	+RKD	94.36	94.29	92.25	86.81	91.91	86.59	

where $C_2, C_3 > 0$.

Combining Eq.20 and Eq.23 completes the proof.

C EXPRIMENTS

Datasets. We take three benchmark datasets: CIFAR-10, CIFAR-100 Krizhevsky (2009) and Tiny ImageNet Le & Yang (2015) and ImageNet Deng et al. (2009). CIFAR-10 has 10 categories while CIFAR-100 has 100 categories. Both datasets consist of 50,000 training samples and 10,000 test samples, with an image resolution of 32×32. Tiny ImageNet contains 100000 images of 200 classes (500 for each class) downsized to 64×64 colored images. Each class has 500 training images, 50 validation images and 50 test images. ImageNet consists of over 1.2 million training images, 50,000 validation images, and 100,000 test images across 1,000 categories with a typical size of 256×256 pixels. Owing to space limitations, the experimental results for CIFAR-10 are provided in the appendix C.

Implementation Details. We set batchsize as 256 and the base learning rate as 0.5 for the teacher model, and set batchsize as 128 and the base learning rate as 0.05 for the student model. We adapt multi-step learning rate decay strategy. Both the teacher and student models are trained for 90 epochs, with the learning rate being reduced three times at epochs [30, 60]. When training the student model, the weight of the teacher's soft labels is set to 0.9. Stochastic Gradient Descent (SGD) is used as the optimizer for experiments. The weight λ was set in the range of 0.01 to 0.03.

Other Results. Next, we present a subset of experimental results. Additional Results on CIFAR-100 and Tiny ImageNet for both homogeneous architecture and heterogeneous architecture are shown in Table 4 and 5. Table 6 shows the results on CIFAR-10, demonstrating the effectiveness of our method. This also demonstrates that the generated soft labels contain more dark knowledge, which is in line with our expectations. Table 7 presents the training time after incorporating intra-class contrastive loss. The introduction of the pipeline effectively reduced the computational cost of intra-class contrastive loss. In total, the training time increased by 10% to 15%. The results of the ablation study on margin loss is shown in Table 8 and 9. We also investigate the sensitivity to the hyperparameter λ , which is shown in Figure 3. T-SNE on CIFAR-100 is illustrated in Figure 2.

Table 7: Training time (per epoch) of the teacher model. It presents the comparison of the training time of the teacher models after incorporating intra-class contrastive loss across different architectures.

Architecture	ResNet50	ResNet34	VGG13	WRN-40-2
Vanilla Teacher	97.36	74.41	79.32	135.11
Ours	108.29	84.94	89.39	152.56
Gap	10.93 (11%)	10.53 (14%)	10.07 (13%)	17.45 (13%)

Table 8: Ablation for Margin Loss on CIFAR-100 and Tiny ImageNet, Homogeneous Architec-

		mager (et) recommended the control						
ture.			CIFAR-100		Tiny ImageNet			
	Teacher	ResNet50	WRN-40-2	VGG13	ResNet50	WRN-40-2	VGG13	
w/o morgin loss	Teacher	72.01	70.93	67.34	53.24	54.77	47.68	
w/o margin loss	Student	ResNet34	WRN-16-2	VGG8	ResNet34	WRN-16-2	VGG8	
		73.65	71.60	68.51	56.26	56.43	49.37	
w/ margin loss	Teacher	ResNet50	WRN-40-2	VGG13	ResNet50	WRN-40-2	VGG13	
		76.57	75.98	74.94	68.89	70.96	66.71	
	Student	ResNet34	WRN-16-2	VGG8	ResNet34	WRN-16-2	VGG8	
	Student	79.10	79.09	74.96	70.22	71.09	65.14	

Table 9: Ablation for Margin Loss on CIFAR-100 and Tiny ImageNet, Heterogeneous Architecture.

			CIFAR-100		Tiny ImageNet			
	Teacher	ResNet50	VGG13	WRN-40-2	ResNet50	VGG13	WRN-40-2	
w/o margin loss	Teacher	72.01	67.34	70.93	53.24	47.68	54.77	
w/o margin loss	Student	MobileNetV2	MobileNetV2	ShuffleV2	MobileNetV2	MobileNetV2	ShuffleV2	
	Student	63.55	60.29	64.83	48.16	45.64	43.18	
	Teacher	ResNet50	VGG13	WRN-40-2	ResNet50	VGG13	WRN-40-2	
w/ morain loss		76.57	74.94	75.98	68.89	66.71	70.96	
w/ margin loss	Student	MobileNetV2	MobileNetV2	ShuffleV2	MobileNetV2	MobileNetV2	ShuffleV2	
	Student	68.72	66.44	72.00	57.31	59.12	61.78	

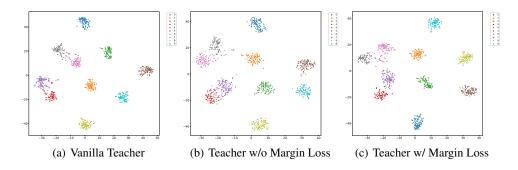


Figure 2: T-SNE on CIFAR-100. We present the t-SNE visualizations of teachers trained with different methods. Compared to the Vanilla Teacher, the incorporation of the intra-class contrastive loss leads to increased intra-class diversity. Furthermore, the comparison between (b) and (c) demonstrates that the margin loss effectively preserves inter-class separation.

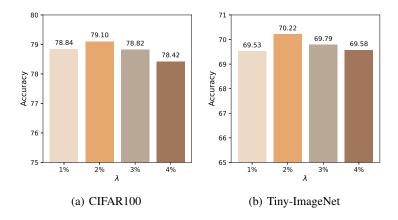


Figure 3: Hyperparameter Sensitivity. Our method performs steadily over different λ . The figures report the accuracy of the student model with homogeneous architecture.