

Correlation Congruence for Knowledge Distillation

Baoyun Peng^{*†1}, Xiao Jin^{*2}, Dongsheng Li¹, Shunfeng Zhou²
Yichao Wu², Jiaheng Liu³, Zhaoning Zhang¹, Yu Liu⁴
¹ NUDT ² SenseTime ³ BUAA ⁴ CUHK
{pengbaoyun13, dsli, zhangzhaoning}@nudt.edu.cn liujiaheng@buaa.edu.cn
{jinxiao, zhoushunfeng, wuyichao}@sensetime.com yuliu@ee.cuhk.edu.hk

Abstract

Most teacher-student frameworks based on knowledge distillation (KD) depend on a strong congruent constraint on instance level. However, they usually ignore the correlation between multiple instances, which is also valuable for knowledge transfer. In this work, we propose a new framework named correlation congruence for knowledge distillation (CCKD), which transfers not only the instance-level information, but also the correlation between instances. Furthermore, a generalized kernel method based on Taylor series expansion is proposed to better capture the correlation between instances. Empirical experiments and ablation studies on image classification tasks (including CIFAR-100, ImageNet-1K) and metric learning tasks (including ReID and Face Recognition) show that the proposed CCKD substantially outperforms the original KD and achieves state-of-the-art accuracy compared with other SOTA KD-based methods. The CCKD can be easily deployed in the majority of the teacher-student framework such as KD and hint-based learning methods. Our code will be released, hoping to nourish our idea to other domains.

1. Introduction

Over the past few decades, various deep neural network (DNN) models have achieved state-of-the-art performance in many vision tasks [36, 37, 14]. Generally, networks with many parameters and computations perform superior to those with fewer parameters and computations when trained on the same dataset. Nevertheless, it's difficult to deploy such large networks on resource-limited embedded systems. Along with the increasing demands for low-cost networks running on embedded systems, there is an urgency for getting smaller network with less computation

and memory consumptions, while narrowing the gap of performance between minor network and large network.

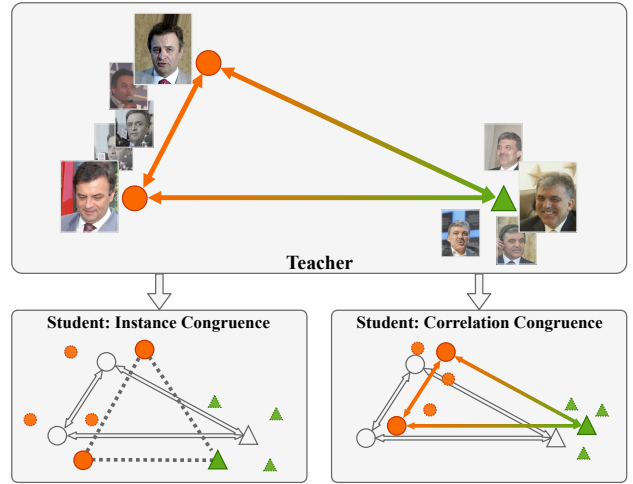


Figure 1: The difference between instance congruence and correlation congruence. When focusing on only instance congruence, the correlation between instances of student may be much different from the teacher's, and the cohesiveness of intra-class would be worse. CCKD solve the problem by adding a correlation congruence when transferring knowledge.

Several techniques have been proposed to address this issue, e.g. parameter pruning and sharing [17, 32], compact convolutional filters [46, 22], low-rank factorization [25, 10] and knowledge distillation [21]. Among these approaches, knowledge distillation has been proved as an effective way to promote the performance of small network by mimicking the behavior of a high-capacity network. It works by adding a strong congruent constraint on outputs of teacher and student for each input instance to encourage the student to mimic teacher's behavior, e.g. minimizing the Kullback-Leibler divergence of predictions [21] or minimizing the euclidean distance of feature representations [29] between teacher and student.

^{*}Equal contribution.

[†]This work was done when Baoyun Peng was an intern at SenseTime Inc.

However, it's hard for the student to learn a mapping function identical to the teacher's due to the gap (in capacity) between teacher and student. By focusing on only instance congruence, the student would learn a much more different instances correlation from the teacher as shown in figure 1. Usually, the embedding space of teacher possesses the characteristic that intra-class instances cohere together while inter-class instances separate from each other. But its counterpart of student model trained by instance congruence would lack such desired characteristic.

We claim that beyond instance congruence, the correlation between instances is also valuable knowledge for promoting the performance of students. Based on this philosophy, we propose a new distillation framework called Correlation Congruence Knowledge Distillation (CCKD) which focus on not only instance congruence, but also correlation congruence to transfer the correlation knowledge between instances to the student as shown in Figure 1. CCKD can be easily implemented and trained with mini-batch, and only requires the same dimension of embedding space for teacher and student network. To cope with the mismatch of feature representations of teacher student network on image classification tasks, we apply a fully-connected layer with the same dimension for both teacher and student network. We conduct various experiments on four representative tasks and different networks to validate the effectiveness of the proposed approach.

Our contributions in this paper are summarized as follows:

1. We propose a new distillation framework named CCKD, which focuses on not only instance congruence but also correlation congruence. To the best of our knowledge, it is the first work to introduce correlation congruence to distillation;
2. We provide a theoretical generalization analysis for CCKD, proving that its feature alignment term reduces the student's RKHS norm and leads to a tighter generalization bound than standard distillation.
3. We introduce a general kernel-based method to better capture the correlation between instances in a mini-batch. We have evaluated and analyzed the impact of different correlation metrics on different tasks;
4. We explore different sampler strategies for mini-batch training to further improve the correlation knowledge transfer;
5. Extensive empirical experiments and ablation studies show the effectiveness of proposed method in different tasks (CIFAR-100, ImageNet-1K, person re-identification and face recognition) to improve the distillation performance.

2. Related Work

Since this paper focuses on training a small but high performance network based on knowledge distillation, we discuss related works in model compression and acceleration, knowledge distillation in this section. In both areas, there are various approaches have been proposed over the past few years. We summarize them as follows.

Model Compression and Acceleration. Model compression and acceleration aim to create network with few computation and parameters cost meanwhile maintaining high performance. A straight way is to design lightweight but powerful network since the original convolution network has many redundant parameters. For example, depth-wise separable convolution is used to replacing standard convolution for building block in [22]. Pointwise group convolution and channel shuffle are proposed to reduce the burden of computation while maintaining high accuracy in [46]. Another way is network pruning which boosts the speed of inference by pruning the neurons or filters with low importance based on certain criteria [17, 32]. In [25, 10], weights were decomposed through low-rank decomposition to save memory cost. Quantization seeks to use low-precision bits to store model's weights or activation outputs [16, 23, 42].

Knowledge Distillation. Transferring knowledge from a large network to a small network is a classical topic and has drawn much attention in recent years. In [21], Hinton *et al.* propose knowledge distillation (KD), in which the student network was trained by the soft output of an ensemble of teacher networks. Comparing to one-hot label, the output from teacher network contains more information about the fine-grained structure among data, consequently helps the student achieve better performance. Since then, there have been works exploring variants of knowledge distillation. In [3], Ba and Caruana show that the performance of a shallower and wider network trained by KD can approximate to deeper ones. Romero *et al.* [33] propose to transfer the knowledge using not only final outputs but also intermediate ones, and add a regressor on intermediate layers to match different size of teacher's and student's outputs. In [45], the authors propose an attention-based method to match the activation-based and gradient-based spatial attention maps. In [44], the flow of solution procedure (FSP), which is generated by computing the Gram matrix of features across layers, was used for knowledge transfer. To improve the robustness of the student, Sau and Balasubramanian [35] perturb the logits of teacher as a regularization.

Different from above offline training methods, several works adopts collaboratively training strategy. Deep mutual learning [47] conducts distillation collaboratively for peer student models by learning from each other. Anil *et al.* [1] further extend this idea by online distillation of multi-

ples networks. In their work, networks are trained in parallel and the knowledge is shared by using distillation loss to accelerate the training process.

Besides, there are several works utilizing adversarial method to modeling knowledge transfer between teacher and student [43, 19, 20]. In [43], they adopt generative adversarial networks combined with distillation to learn the loss function to better transfer teacher’s knowledge to student. Byeongho *et al.* [20] adopt adversarial method to discover adversarial samples supporting decision boundary.

In this paper, beyond instance knowledge, we take the correlation in embedded space between instances as valuable knowledge to transfer correlation among instances in the embedded space between for knowledge distillation.

3. Correlation Congruence Knowledge Distillation

In this section, we describe the details of proposed method based on correlation congruence for knowledge distillation.

3.1. Background and Notations

We refer a well-performed teacher network with parameters \mathbf{W}_t as T and a new student network with parameters \mathbf{W}_s as S like in [21, 45, 44, 1, 33]. The input dataset of the network is noted as $\chi = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, and the corresponding ground truth is noted as $\mathcal{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$, n represents the number of samples in dataset. Since deep network can be viewed as a mapping function stacked by multiple non-linear layers, we note $\phi_t(\mathbf{x}; \mathbf{W}_t)$ and $\phi_s(\mathbf{x}; \mathbf{W}_s)$ as the mapping functions of teacher and student, \mathbf{x} represents the input data. \mathbf{f}_s and \mathbf{f}_t represent the feature representations of teacher and student. The logits of teacher and student note as $\mathbf{z}_t = \phi(\mathbf{x}; \mathbf{W}_s)$ and $\mathbf{z}_s = \phi(\mathbf{x}; \mathbf{W}_t)$. $\mathbf{p}_t = \text{softmax}(\mathbf{z}_t)$ and $\mathbf{p}_s = \text{softmax}(\mathbf{z}_s)$ represent the final prediction probabilities of teacher and student.

3.2. Knowledge Distillation

Overparameterized networks have shown powerful optimization properties to learn the desired mapping function from data [12], of which the output reflects fine-grained structure one-hot labels might ignore. Based on this insight, knowledge distillation was first proposed in [7] for model compression, then Hinton *et al.* [21] popularized it. The idea of knowledge distillation is to let the student mimic the teacher’s behavior by adding a strong congruent constraint on predictions [7, 21, 33] using KL divergence

$$L_{\text{KD}} = \frac{1}{n} \sum_{i=1}^n \tau^2 \text{KL}(\mathbf{p}_s^\tau, \mathbf{p}_t^\tau), \quad (1)$$

where τ is a relaxation hyperparameter (referred as temperature in [21]) to soften the output of teacher network,

$\mathbf{p}^\tau = \text{softmax}(\frac{\mathbf{z}}{\tau})$. In several works [38, 29] the KL divergence is replaced by euclidean distance,

$$L_{\text{mimic}} = \frac{1}{n} \sum_{i=1}^n \|\mathbf{f}_s - \mathbf{f}_t\|_2^2. \quad (2)$$

Regardless of congruent constraint on final predictions [21], feature representations [38] or activations of hidden layer [33], these methods only focus on instance congruence while ignore the correlation between instances. Due to the gap (in capacity) between teacher and student, it’s hard for student to learn a identical mapping function from teacher by instance congruence. We argue that the correlation between instances is also vital for classification since it directly reflect how the teacher model the structure of different instances in embedded feature space.

3.3. Correlation Congruence

In this section, we describe CCKD in detail. Different from previous methods, CCKD considers not only the instance level congruence but also correlation congruence between instances. Figure 2 shows the overview of CCKD. CCKD consists of two part: instance congruence (KL divergence on predictions of teacher and student) and correlation congruence (euclidean distance on correlation of teacher and student).

Let \mathbf{F}_t and \mathbf{F}_s represent the set of feature representations of teacher and student respectively,

$$\begin{aligned} \mathbf{F}_t &= \text{matrix}(\mathbf{f}_1^t, \mathbf{f}_2^t, \dots, \mathbf{f}_n^t), \\ \mathbf{F}_s &= \text{matrix}(\mathbf{f}_1^s, \mathbf{f}_2^s, \dots, \mathbf{f}_n^s). \end{aligned} \quad (3)$$

The feature \mathbf{f} can be seen as a point in the embedded feature space. Without loss of generality, a mapping function is introduced as follow:

$$\psi : \mathbf{F} \rightarrow \mathbf{C} \in \mathbb{R}^{n \times n}. \quad (4)$$

where \mathbf{C} is a correlation matrix. Each element in \mathbf{C} represents the correlation between \mathbf{x}_i and \mathbf{x}_j in embedding space, which is defined as

$$\mathbf{C}_{ij} = \varphi(\mathbf{f}_i, \mathbf{f}_j), \quad \mathbf{C}_{ij} \in \mathbb{R} \quad (5)$$

The function φ can be any correlation metric, and we will introduce three metric for capturing the correlation between instances in next section. Then, the correlation congruence can be formulated as follow:

$$\begin{aligned} L_{\text{CC}} &= \frac{1}{n^2} \|\psi(\mathbf{F}_t) - \psi(\mathbf{F}_s)\|_2^2 \\ &= \frac{1}{n^2} \sum_{i,j} (\varphi(\mathbf{f}_i^s, \mathbf{f}_j^s) - \varphi(\mathbf{f}_i^t, \mathbf{f}_j^t))^2. \end{aligned} \quad (6)$$

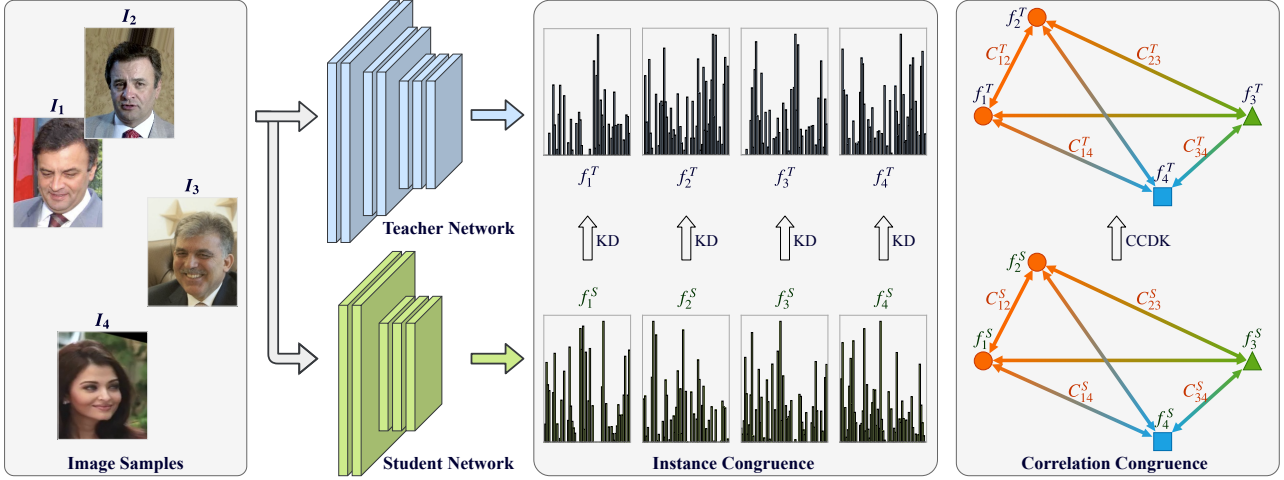


Figure 2: The overall framework of correlation congruence for knowledge distillation (T : teacher; S : student; f_i^T : teacher’s output of i_{th} sample; f_i^S : student’s output of i_{th} sample; C_i : correlation between i_{th} and j_{th} sample). Original KD focus on only instance congruence between teacher and student network. While CCKD aims to not only instance congruence but also correlation congruence between multiple instances.

Then, the optimization goal of CCKD is to minimize the following loss function:

$$L_{CCKD} = \alpha L_{CE} + (1 - \alpha) L_{KD} + \beta L_{CC}, \quad (7)$$

where L_{CE} is the cross-entropy loss, α and β are two hyper-parameters for balancing correlation congruence and instance correlation.

3.4. Generalized kernel-based correlation

Capturing the complex correlations between instances is not easy due to a very high dimension in the embedded space [39]. In this section, we introduce kernel trick to capture the high order correlation between instances in the feature space.

Let $\mathbf{x}, \mathbf{y} \in \Omega$ represent two instances in feature space, and we introduce different mapping functions $k : \Omega \times \Omega \mapsto \mathbb{R}$ as correlation metric, including:

1. naive MMD: $k(\mathbf{x}, \mathbf{y}) = \left| \frac{1}{n} \sum_i \mathbf{x}_i - \frac{1}{n} \sum_i \mathbf{y}_i \right|$;
2. Bilinear Pool: $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}$;
3. Gaussian RBF: $k(\mathbf{x}, \mathbf{y}) = \exp(-\frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{2\delta^2})$;

MMD can reflect the distance between mean embeddings. Bilinear Pooling [30] can be seen as a naive 2-th order function, of which the correlation between two instances is computed by element-wise dot product. Gaussian RBF is a common kernel function whose value depends only on the euclidean distance from the origin space.

Comparing to naive MMD and Bilinear Pool, Gaussian RBF is more flexible and powerful in capturing the complex non-linear relationship between instances. Based on Gaussian RBF, the correlation mapping function ϕ can be computed by a kernel function $K : F \times F \rightarrow \mathbb{R}^{n \times n}$, where each element can be computed as

$$[k(\mathbf{F}, \mathbf{F})]_{ij} \approx \sum_{p=0}^P \alpha_p (\mathbf{F}_i \cdot \mathbf{F}_{j\cdot}^\top)^P. \quad (8)$$

which can be approximated by P -order Taylor series. Once specifying the kernel function, then the coefficient α_p is also confirmed. Each element $[k(\mathbf{F}, \mathbf{F})]_{ij}$ encodes the pairwise correlations between i -th and j -th features in \mathbf{F} . We take Gaussian RBF kernel function as an example, then

$$\begin{aligned} [k(\mathbf{F}, \mathbf{F})]_{ij} &= \exp(-\gamma \|\mathbf{F}_i - \mathbf{F}_{j\cdot}\|^2) \\ &\approx \sum_{p=0}^P \exp(-2\gamma) \frac{(2\gamma)^p}{p!} (\mathbf{F}_i \cdot \mathbf{F}_{j\cdot}^\top)^p. \end{aligned} \quad (9)$$

where γ is a tunable parameter.

3.5. Strategy for Mini-batch Sampler

Usually, stochastic gradient descent (SGD), which samples batch of training examples uniformly at random from training dataset, is adopted to train the network and then parameters are updated using the sampled batch of examples. The naive random sampler would lead to such a situation that all examples come from different classes. Although it is an unbiased estimation to true gradient of instance congruence, uniformly sampled mini-batch gradient descent will

result in a high biased estimation to gradient of intra-class correlation. To transfer the true correlation information unbiasedly, a proper sampler strategy is important.

To balance the intra-class and inter-class correlation congruence, we propose two strategies for mini-batch sampler: class-uniform random sampler (CUR-sampler) and superclass-uniform random sampler (SUR-sampler). CUR-sampler samples by class and randomly selects fixed k number of examples for each sampled class (eg. each batch consists of 6 class and each class contains $k = 8$ examples, forming a 48 batch size). SUR-sampler is similar to CUR-sampler, but different in that it samples examples by superclass, a more soft form of true class generated by clustering. To get the superclass of training examples, we first extract the feature using teacher model, then use the K-means to cluster. The superclass of example is defined as the cluster it belongs to. Comparing to CUR-sampler, SUR-sampler is more flexible and tolerant for imbalance label since the superclass reflects the coarse structure of instances in embedded space.

3.6. Complexity analysis and details

To cope with the mini-batch training, we compute the correlation in a mini-batch. Formula 9 involves the computation of a large pairwise matrix $b \times b$ (b is the batch size), and each element is approximated by p -order Taylor-series with p times dot product computation between two d dimension vectors. The total computation complexity is $O(pbd^2)$ in a mini-batch, and the extra space consumption is $O(b^2 + d^2)$ for storing the correlation matrix. Compared to huge parameters and computation for training deep neural network, the time and computation consumption for correlation congruence can be ignored. Besides, since the correlation congruent constraint is added on embedding space, it only requires that the feature dimension of student network is the same as teacher. To cope with the mismatch dimension in classification tasks, a fully-connected layer with fixed-length dimension is added for both teacher and student network, which has minor influence on other methods in this paper.

4. Experiments

We evaluate CCKD on multiple tasks, including image classification tasks (CIFAR-100 and ImageNet-1K) and metric learning tasks (including MSMT17 dataset ReID and MegaFace for face recognition), and compare it with closely related works. Extensive experiments and analysis are conducted to delve into the CCKD.

4.1. Experimental Settings

Network Architecture and Implementation Details
Given the steady performance and efficiency computation,

ResNet [18] and MobileNet [34] network are chosen in this work.

In the main experiments, we set the order $P = 2$, and compute Equation 9 in a mini batch. For the networks in CIFAR-100 and ImageNet-1K, we add a fully-connected layer with 128-d output to form a sharing embedding space for teacher and student. The hyper-parameter α is set to zero, and correlation congruence scale β is set to 0.003, $\gamma = 0.4$. CUR-sampler is used for all the main experiments with $k = 4$.

On CIFAR-100, ImageNet-1K and MSMT17, Original Knowledge distillation (KD) [21] and cross-entropy (CE) are chosen as the baselines. For face recognition, ArcFace loss [9] and $L2$ -mimic loss [29, 31] are adopt. We compare CCKD with several state-of-the-art distillation related methods, including attention transfer (AT) [45], deal mutual learning (DML) [47] and conditional adversarial network (Adv) [43]. For attention transfer, we add it for last two blocks as suggested in [45]. For adversarial training, the discriminator consists of FC(128×64) + BN + ReLU + FC (64×2) + Sigmoid activation layers, and we adopt BinaryCrossEntropy loss to train it. All the networks and training procedures are implemented in PyTorch.

4.2. Classification Results on CIFAR-100

CIFAR-100 [27] consists of colored natural images with 32×32 size. There are 100 classes in CIFAR-100, each class contains 500 images in training set and 100 images in validation set. We use a standard data augmentation scheme (flip/padding/random crop) that is widely used for these dataset, and normalize the input images using the channel means and standard deviations. We set the weight decay of student network to $5e-4$, batch size to 64, and use stochastic gradient descent with momentum. The starting learning rate is set as 0.1, and divided by 10 at 80, 120, 160 epochs, totally 200 epochs. Top-1 and top-5 accuracy are adopted as performance metric.

Table 1: Validation accuracy results on CIFAR-100. ResNet-110 is as teacher network, ResNet-20 and ResNet-14 as student networks. We keep the same training configuration for all the methods for fair comparison.

method	resnet-20		resnet-14	
	top-1	top-5	top-1	top-5
CE	68.4	91.3	66.4	90.3
KD	70.8	92.4	68.3	90.7
DML	71.2	92.5	69.1	91.2
AT	71.0	92.4	68.6	91.1
Adv	70.5	92.1	68.1	90.6
CCKD	72.4	92.9	70.2	92.0

Table 1 summarizes the results of CIFAR-100. CCKD

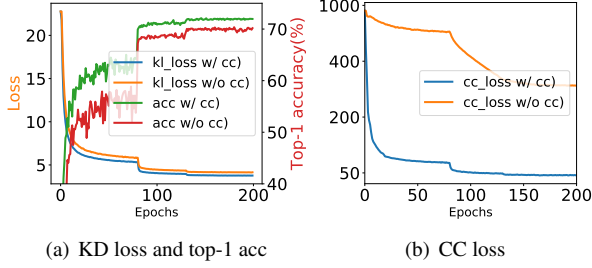


Figure 3: The curve of training loss and validation accuracy.

gets a 72.4% and 70.2% of top-1 accuracy for ResNet-20 and ResNet-14, and substantially surpasses the CE by 4.0% and 3.8%, 1.6% and 1.9% over KD. For the online distillation DML [47], we train target network (ResNet-14 and ResNet-20) collaboratively with ResNet-110, and evaluate performance of target network. Comparing to other SOTA methods, CCKD still significantly outperforms them. All the four distillation related methods significantly surpass the original CE over 2%, which verifies the effectiveness of teacher-student methods.

Figure 3 shows the training loss and validation accuracy of ResNet-20. It can be observed that although KL divergence loss after convergence is almost the same, the correlation congruence loss for CCKD is much lower than original KD, consequently results in a higher performance.

4.3. Results on ImageNet-1K

ImageNet-1K [8] consists 1.28M training images and 50K testing images in total. We adopt the ResNet-50 [18] as the teacher network, MobileNetV2 with 0.5 width multiplier as the student network. The data augmentation scheme for training images is the same as [18], and apply a center-crop at test time. All the images are normalized using the channel means and standard deviations. We set the weight decay of student network to $5e-4$, batch size to 1,024 (training on 16 TiTAN X, each with 64 batch size), and use stochastic gradient descent with momentum. The starting learning rate is set as 0.4, then divided by 10 at 50, 80, 120 epochs, totally 150 epochs.

For fair comparison, we keep the same configuration for all the methods. Table 2 summarizes the results on ImageNet 1K. CCKD gets a 67.7% Top-1 accuracy, which surpasses the cross-entropy by 3.3%. Compare with original KD[21], CCKD surpasses by 1.0% in top-1 accuracy. AT and DML perform worse than original KD. To our best knowledge, we have not found any works that successfully verify the effectiveness of KD on ImageNet-1K dataset. It has been reported in work [45] that KD struggles to work when the architecture and depth of student network are different from the teacher. But we found that by removing the dropout layer and using a proper temperature (T in [4,8]), the KD can surpass the student over 2.0%.

Table 2: Validation accuracy results on ImageNet 1K. The teacher network is ResNet-50, student network is MobileNetV2 with 0.5 width multiplier. We keep the same configuration for CE and other four student networks.

method	top-1 accuracy	top-5 accuracy
teacher	75.5	92.7
CE	64.2	85.4
KD	66.7	87.3
DML	65.3	86.1
Adv	66.8	87.3
AT	65.4	86.1
CCKD	67.7	87.7

4.4. Person Re-Identification on MSMT17

Comparing to closed set classification, open set classification is more dependent on a good metric learning and more realistic scenario. We apply the proposed method to two open-set classification: person re-identification (ReID) and face recognition.

For ReID, we evaluate proposed method on MSMT17 [41]. It contains 180 hours of videos captured by 12 outdoor cameras, 3 indoor cameras under different seasons and time. There are 126,441 bounding boxes of 4,101 identities that are annotated. All the bounding boxes are split to training set (32621 bounding boxes, 1041 identities), query set (11659 bounding boxes, 3060 identities) and gallery set (82161 bounding boxes). There is no intersection of identities between training set and query & gallery set. We train the networks on training set, and perform identification on query and gallery set. Rank-1&5 and mean average precision (mAP) are adopted as performance metric.

ResNet-50 is used as the teacher network and ResNet-18 as student network. The dimension of the feature representation is set to 256. We set the weight decay to $5e-4$, batch size to 40, and use stochastic gradient descent with momentum. The learning rate is set as 0.0003, then divided by 10 at 45, 60 epochs, totally 90 epochs.

Table 3: Validation accuracy results on MSMT17. The teacher network is ResNet-50, student network is Resnet-18.

method	pretrained?	rank-1	rank-5	mAP
teacher	yes	66.4	79	34.3
CE	no	32.4	49.0	14.2
DML-1	no	34.5	51.5	16.5
DML-2	yes	50.2	66.4	25.3
KD	no	56.8	72.3	28.3
AT	no	57.6	72.5	28.7
Adv	no	56.0	71.6	27.8
CCKD	no	59.7	74.1	30.7

Table 3 summarizes the results of MSMT17 with CCKD,

as well as the comparison against other SOTA methods. For fair comparison, all the distillation based methods (except DML) are trained without ImageNet-1K pretraining. For DML, both the results with/without ImageNet-1K pretraining are represented. It can be seen that the performance of the CCKD significantly surpasses KD and other SOTA KD-based methods, and promotes the original KD by 3.1% for rank-1 accuracy and 2.4% for mAP. Without the guidance of teacher, the student trained by cross-entropy only achieves 14.2% mAP, which is much lower than 28.3% of KD.

4.5. Face recognition results on Megaface

Similar to ReID, face recognition is a classical metric learning problem. Learning a discriminative embedded space is the key to get a powerful recognition model. Usually, thousands of identities (class) are required for training a well-performed recognition model. Empirical evidence shows that mimicking the feature layer with hint-based L2 Loss can bring great improvement for small network [29, 31]. In this experiment, instead of using KD loss, we adopt the L2-mimic loss. MS-Celeb-1M [15] and IMDB-Face [40] are used as training datasets.

We choose MegaFace [26], a very popular benchmark, as testing set to evaluate the proposed method. MegaFace aims at the evaluation of face recognition algorithms at million-scale of distractors (people who are not in the testing set). We adopt 1:N identification protocol in Megaface to evaluate the different methods. Rank-1 identification rate at different number of distractors is used as metric for evaluation. We set weight decay to $5e-4$, batch size to 1024, and use stochastic gradient descent with momentum. The learning rate is set as 0.1, and divided by 10 at 50, 80, 100 epochs, 120 epochs in total. ResNet-50 is used as teacher network, and MobileNetV2 with 0.5 width multiplier as student network.

Table 4: Results on Megaface. The teacher network is ResNet-50 trained on MsCeleb-1M [15] and IMDB-face [40] using ArcFace [9]. The student network is MobileNetV2 with a width multiplier=0.5. We keep the same training configuration for mimic, mimic with Adv and CCKD.

method	Rank-1 Identification rate at different distractors					
	ds=10 ¹	ds=10 ²	ds=10 ³	ds=10 ⁴	ds=10 ⁵	ds=10 ⁶
teacher	99.76	99.66	99.58	99.49	99.23	98.15
student	99.20	96.37	91.49	84.45	75.60	65.91
mimic	99.63	98.73	97.25	94.39	89.60	83.01
mimic+Adv	99.64	98.80	97.43	94.81	90.52	84.13
CCKD	99.66	99.07	97.93	95.76	91.99	86.29

Table 4 shows the results on Megaface. It can be observed that ArcFace loss, which is trained by only using pure one-hot labels, achieves 65.91% Rank-1 identification rate with 1M distractors. When guided by the teacher using L2-mimic loss, the student network can achieves 83.01%,

promoting by 18.1%. This result shows that even a much small network can get a substantial improvement of performance when designing proper target and optimization goal. By adding the constraints on correlations among instance, CCKD achieves 86.29% Rank-1 identification rate with 1M distractors, which surpasses the mimicking by 3.28% and 2.16% promotion over Adv [43].

4.6. Ablation Studies

Correlation Metrics. To explore the impact of different correlation metrics on CCKD, we evaluate three popular metrics, namely max mean discrepancy (MMD), Bilinear Pool and Gaussian RBF. We approximate the Gaussian RBF by using 2-order Taylor series. MMD reflects the difference between two instances in mean embeddings. Bilinear Pool evaluate the similarity of instances pair, and we adopt identity matrix as the linear matrix. When the features are normalized to unit length, it is equal to the cosine similarity. Gaussian RBF is a common kernel function whose value depends only on the euclidean distance from the origin space.

Table 5: Results on MSMT17 with different correlation methods, including MMD, Bilinear Pool and Gaussian RBF. The Gaussian RBF achieves the best result.

correlation metric	rank-1	rank-5	mAP
MMD	58.9	73.6	29.4
Bilinear	59.2	73.8	30.2
Gaussian RBF	59.6	74.0	30.4

Table 5 shows the results of MSMT17 with different correlation metrics. Gaussian RBF achieves the better performance comparing to MMD and Bilinear Pool, while MMD performs worst. So in the main experiments, we use the Gaussian RBF approximated by 2-order Taylor series. All the three correlation matrices greatly surpass the original KD, which proves the effectiveness of correlation in knowledge distillation.

Order of Taylor series. To exploit the high order of correlations between instances, we expand the Gaussian RBF by Taylor series to 1, 2, 3 -order respectively.

Table 6: Results on MSMT17 with different order ($p = 1, 2, 3$) Taylor series.

Expand order	rank-1	rank-5	mAP
$p=1$	59.2	73.7	30.1
$p=2$	59.6	74	30.4
$p=3$	60.5	74.5	30.7

Table 6 summarizes the results on MSMT17 with ap-

proximated Gaussian RBF at different orders. It can be observed that 3-order is better than 1, 2-order, and 1-order performs worst. Generally speaking, expanding Gaussian RBF to high order can capture more complex correlations, and consequently achieves higher performance in knowledge distillation.

Impact of Different Sampler Strategies. To explore a proper sampler strategy, we evaluate the impacts of different sampler strategies including uniform random sampler (UR-sampler), class-uniform random sampler (CUR-sampler) and superclass-uniform random sampler (SUR-sampler) on MSMT17 dataset. For SUR-sampler, the k-means is adopted and the number of clusters is set to 1000 to generate superclass. For fair comparison the batchsize is set to 40 for all three strategies, and we set different $k = 1, 2, 4, 8, 20$ both for CUR-sampler and SUR-sampler.

Table 7: Results on MSMT17 with different batch sampler strategies. The teacher network is ResNet-50 and the student network is ResNet-18.

sampler	rank-1	rank-5	mAP
UR-sampler	57.2	72.3	28.6
CUR-sampler($k=1$)	57.4	72.4	28.8
CUR-sampler($k=2$)	58.9	73.6	29.4
CUR-sampler($k=4$)	59.7	74.1	30.2
CUR-sampler($k=8$)	55.7	71.8	29.1
CUR-sampler($k=20$)	24.7	40.9	10.7
SUR-sampler($k=1$)	56.2	72.2	29.4
SUR-sampler($k=2$)	58.3	73.9	29.9
SUR-sampler($k=4$)	59.6	75.0	31.1
SUR-sampler($k=8$)	56.2	72.2	29.4
SUR-sampler($k=20$)	30.1	47.7	13.7

Table 7 summarizes the results. It can be observed that the sampler strategy have a great impact on performance. Both SUR-sampler and CUR-sampler are sensitive to the value of k , which plays a role of balancing the intra-class and inter-class correlation congruence. When given fixed batch size, a larger k means a smaller number of classes in a mini-batch. Both CUR-sampler and SUR-sampler become worse when $k = 8$ or above. A possible explanation is that small number of classes in a mini-batch results a high bias estimation for true gradient. While the SUR-sampler performs better than CUR-sampler in such bad cases. By selecting proper k (eg. 2 or 4 in our experiments), Both CUR-sampler and SUR-sampler performs better than UR-sampler.

4.7. Analyze

To delving into essence beyond results, we perform analysis based on visualization. We count the cosine similarities of intra-class instances and inter-class instances on

MSMT17 since it is a common metric for openset recognition. Figure 4 shows the heatmaps of cosine similarities. The top row shows intra-class instances and the bottom row shows inter-class instances from two different identities. Each cell relates to cosine similarity between corresponding instance pair.

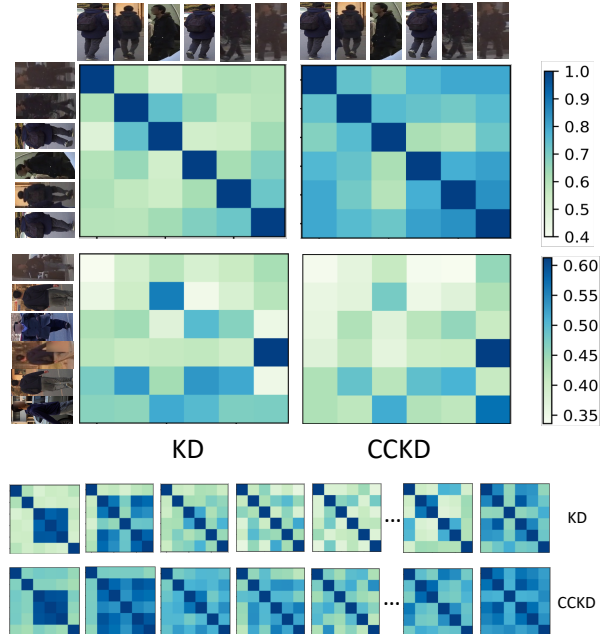


Figure 4: The heatmaps of cosine similarities between instances pairs. The top row shows intra-class similarities and the middle row shows inter-class similarities between two identities. More intra-class heatmap are showed in bottom two rows. (best viewed in color)

It can be observed that, cosine similarity between intra-class instances of CCKD is more larger than KD overall, which means a more cohesion of intra-class instances in embedding space. Although there is not much difference between CCKD and KD in inter-class cosine similarity. It seems that CCKD can help the student to learn a more discriminative embedding space. While CCKD by considering the correlation congruence between instances, consequently getting a better performance.

5. Conclusions

In this paper, we propose a new distillation framework named CCKD, which considers not only instance information but also correlation information between instances when transferring knowledge. To better capture correlation, a generalized method based on Taylor series expansion of kernel function is proposed. To further improve the CCKD, two new mini-batch sampler strategies are proposed. Extensive experiments on four representative tasks show that

the proposed approach can significantly promote the performance of student network.

Appendix: Generalization Analysis for CCKD

Here, we provide a theoretical generalization analysis for CCKD. The proof is built on the Neural Tangent Kernel (NTK) framework. More specifically, we use stability analysis of optimization dynamics and derive a new geometric condition, and under this condition, CCKD reduces the norm in the RKHS. This lets us prove a formal generalization bound. It shows how CCKD’s feature alignment improves generalization over standard KD [21]. We clarify the link between feature matching, loss curvature, and model complexity.

A.1 Preliminaries and Notation

We define the math setup and notation here.

Definition 1 (Setup). Let $\mathcal{X} \subseteq \mathbb{R}^p$ be the input space, $\mathcal{Y} = \{1, \dots, K\}$ be the label space. Data comes from an unknown distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$. Training set $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathcal{D}^n$. A neural network is a function $h : \mathcal{X} \rightarrow \mathbb{R}^K$. We study two networks: a fixed teacher h_T and a student h_S .

Definition 2 (Neural Tangent Kernel (NTK) and RKHS). For student networks parameterized by Φ , the NTK is:

$$\Theta(\mathbf{x}, \mathbf{x}') = \text{tr}(\nabla_{\Phi} h(\mathbf{x})^{\top} \nabla_{\Phi} h(\mathbf{x}'))$$

It defines a unique RKHS \mathcal{H}_{Θ} . This space has norm $\|\cdot\|_{\mathcal{H}_{\Theta}}$ and inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}_{\Theta}}$ [24, 2].

Definition 3 (Feature Gram Matrix). Let $f : \mathcal{X} \rightarrow \mathbb{R}^d$ be the feature extractor (before the final classifier). The Gram matrix on S is $[\mathbf{K}_S(h)]_{ij} = \langle f(\mathbf{x}_i), f(\mathbf{x}_j) \rangle \in \mathbb{R}^{n \times n}$. We write $\mathbf{K}_S(h_S)$ for the student. We write \mathbf{K}_T for the pre-computed teacher Gram matrix.

Definition 4 (Learning Objectives and Risk). Let ℓ_{CE} be the Cross-Entropy loss, ℓ_{KL} be the KL divergence for soft targets, $\alpha \in [0, 1]$ be the distillation weight, $\mu > 0$ be the weight decay parameter.

1. **KD** mixes classification and distillation losses. It obtains h_S^{KD} by minimizing:

$$h_S^{KD} := \arg \min_{h \in \mathcal{H}_{\Theta}} \left\{ \hat{R}_n(h) + \frac{\mu}{2} \|h\|_{\mathcal{H}_{\Theta}}^2 \right\}, \quad (10)$$

$$\text{where } \hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \left[(1 - \alpha) \ell_{CE}(h(\mathbf{x}_i), y_i) + \alpha \ell_{KL}(h(\mathbf{x}_i), h_T(\mathbf{x}_i)) \right] \quad (11)$$

2. **CCKD** incorporates a feature alignment term. It penalises the Frobenius norm between the student and teacher Gram matrices, and finds h_S^{CCKD} by minimizing:

$$h_S^{CCKD} := \arg \min_{h \in \mathcal{H}_{\Theta}} \left\{ \hat{R}_n(h) + \frac{\mu}{2} \|h\|_{\mathcal{H}_{\Theta}}^2 + \underbrace{\frac{\lambda}{n^2} \|\mathbf{K}_S(h) - \mathbf{K}_T\|_F^2}_{\text{Feature Alignment Loss } G(h)} \right\} \quad (12)$$

Here, $\lambda > 0$ controls the strength of the alignment.

For simplicity, we introduce the following shorthands:

- $J_{KD}(h) = \hat{R}_n(h) + \frac{\mu}{2} \|h\|^2$
- $J_{CCKD}(h) = J_{KD}(h) + \lambda G(h)$

A.2 Assumptions

Assumption 1 (Kernel Regime). The training process remains within the NTK regime, where the kernel Θ remains fixed. The function h_S evolves inside \mathcal{H}_{Θ_S} [28].

Assumption 2 (Smoothness and Local Convexity). Under Assumption 1, the loss is convex in network outputs. We assume $\hat{R}_n(h)$ is convex and twice differentiable. We assume $G(h)$ is convex and twice differentiable, and similarly $G(h)$ is **convex near** h_S^{KD} .

Remark 1 (On Local vs. Global Minimizers). Local convexity of $G(h)$ is a common assumption when analysing complex non-convex objectives. It implies $J_{CCKD}(h)$ is strongly convex near h_S^{KD} . By the implicit function theorem [11], a unique local minimizer $h(\lambda)$ exists near $h(0) = h_S^{KD}$. We call this h_S^{CCKD} .

Our method tracks the solution path $h(\lambda)$ as λ increases from zero, analogous to the “regularization path” analysis in classical ML, such as LASSO [13]. We focus on the minimiser identified along this path, which corresponds to the typical practical scenario when λ is small.

Analysing local curvature is a standard technique in deep learning theory, particularly within the NTK framework [5]. For sufficiently small λ , both h_S^{CCKD} and h_S^{KD} lie in the same locally convex region. Therefore, $J_{CCKD}(h_S^{CCKD}) \leq J_{CCKD}(h_S^{KD})$ holds, which supports the subsequent proofs.

Assumption 3 (Hessian-Transformed Geometric Alignment). Let $H_{J_{KD}}(h) = H_{\hat{R}_n}(h) + \mu I$ be the KD Hessian. At h_S^{KD} , we assume:

$$\langle h_S^{KD}, (H_{J_{KD}}(h_S^{KD}))^{-1} \nabla_h G(h_S^{KD}) \rangle_{\mathcal{H}_{\Theta_S}} > 0 \quad (13)$$

Remark 2 (Interpretation of Assumption 3). This condition implies that CCKD effectively reduces the norm along directions where the KD loss landscape is relatively flat or “soft”.

Assumption 4 (Lipschitz Loss). *The empirical risk $\hat{R}_n(h)$ is L_R -Lipschitz with respect to the function outputs over the training set [4].*

A.3 Main Results

Lemma 1 (CCKD-Induced Norm Reduction). *Under Assumptions 1–3, there exists $\lambda_0 > 0$. For all $0 < \lambda < \lambda_0$, the CCKD solution satisfies:*

$$\|h_S^{\text{CCKD}}\|_{\mathcal{H}_{\Theta_S}} < \|h_S^{\text{KD}}\|_{\mathcal{H}_{\Theta_S}} \quad (14)$$

Detailed Proof. Step 1: Setup for Perturbation Analysis. Let $h(\lambda) = h_S^{\text{CCKD}}$ minimize $J_{\text{CCKD}}(h)$ for $\lambda \geq 0$. Note: $h(0) = h_S^{\text{KD}}$. The function $h(\lambda)$ satisfies:

$$\nabla_h \hat{R}_n(h(\lambda)) + \mu h(\lambda) + \lambda \nabla_h G(h(\lambda)) = 0 \quad (15)$$

We want to see how $\|h(\lambda)\|^2$ changes as λ grows from 0. We analyse the sign of its derivative at $\lambda = 0$.

Step 2: Differentiating the Squared Norm. Let $\phi(\lambda) = \|h(\lambda)\|^2$. Then:

$$\phi'(\lambda) = 2 \langle h(\lambda), h'(\lambda) \rangle \quad (16)$$

Step 3: Finding $h'(\lambda)$ via Implicit Differentiation. Differentiate Eq. 15 w.r.t. λ . Apply the chain rule and product rule:

- $H_{\hat{R}_n}(h(\lambda))h'(\lambda)$
- $\mu h'(\lambda)$
- $\nabla_h G(h(\lambda)) + \lambda H_G(h(\lambda))h'(\lambda)$

Combine them:

$$(H_{\hat{R}_n}(h(\lambda)) + \mu I + \lambda H_G(h(\lambda)))h'(\lambda) = -\nabla_h G(h(\lambda)) \quad (17)$$

Step 4: Evaluating at $\lambda = 0$. At $\lambda = 0$, $h(0) = h_S^{\text{KD}}$. So:

$$(H_{\hat{R}_n}(h_S^{\text{KD}}) + \mu I)h'(0) = -\nabla_h G(h_S^{\text{KD}}) \quad (18)$$

By definition, $H_{J_{\text{KD}}}(h_S^{\text{KD}}) = H_{\hat{R}_n}(h_S^{\text{KD}}) + \mu I$. This Hessian is invertible (strong convexity) [6]. Solve for $h'(0)$:

$$h'(0) = -(H_{J_{\text{KD}}}(h_S^{\text{KD}}))^{-1} \nabla_h G(h_S^{\text{KD}}) \quad (19)$$

Step 5: Sign of the Norm's Derivative. Plug Eq. 19 into Eq. 16:

$$\begin{aligned} \phi'(0) &= 2 \langle h_S^{\text{KD}}, h'(0) \rangle \\ &= -2 \langle h_S^{\text{KD}}, (H_{J_{\text{KD}}}(h_S^{\text{KD}}))^{-1} \nabla_h G(h_S^{\text{KD}}) \rangle_{\mathcal{H}_{\Theta_S}} \end{aligned} \quad (20)$$

Step 6: Conclude with Assumption 3. Assumption 3 says the inner product is positive. So $\phi'(0) < 0$.

By Assumption 2, all functions are smooth. Then $h(\lambda)$ is continuously differentiable. So is $\phi(\lambda)$.

Since $\phi'(0) < 0$, there exists $\lambda_0 > 0$. For $\lambda \in (0, \lambda_0)$, we have $\phi(\lambda) < \phi(0)$. Thus:

$$\|h_S^{\text{CCKD}}\|^2 < \|h_S^{\text{KD}}\|^2$$

□

Lemma 2 (Feature Alignment Guarantee). *For any $\lambda > 0$, $G(h_S^{\text{CCKD}}) \leq G(h_S^{\text{KD}})$. If $\nabla_h G(h_S^{\text{KD}}) \neq 0$, the inequality is strict for small $\lambda > 0$.*

Detailed Proof. Part 1: Non-strict Inequality.

By definition:

$$J_{\text{KD}}(h_S^{\text{KD}}) \leq J_{\text{KD}}(h_S^{\text{CCKD}}) \quad (21)$$

$$J_{\text{CCKD}}(h_S^{\text{CCKD}}) \leq J_{\text{CCKD}}(h_S^{\text{KD}}) \quad (22)$$

Expand Eq. 22:

$$J_{\text{KD}}(h_S^{\text{CCKD}}) + \lambda G(h_S^{\text{CCKD}}) \leq J_{\text{KD}}(h_S^{\text{KD}}) + \lambda G(h_S^{\text{KD}})$$

Use Eq. 21:

$$\begin{aligned} J_{\text{KD}}(h_S^{\text{KD}}) + \lambda G(h_S^{\text{CCKD}}) &\leq J_{\text{KD}}(h_S^{\text{CCKD}}) + \lambda G(h_S^{\text{CCKD}}) \\ &\leq J_{\text{KD}}(h_S^{\text{KD}}) + \lambda G(h_S^{\text{KD}}) \end{aligned} \quad (23)$$

Simplify: $\lambda G(h_S^{\text{CCKD}}) \leq \lambda G(h_S^{\text{KD}})$.

Since $\lambda > 0$, we get: $G(h_S^{\text{CCKD}}) \leq G(h_S^{\text{KD}})$.

Part 2: Strict Inequality.

Let $\psi(\lambda) = G(h(\lambda))$. Its derivative at $\lambda = 0$ is:

$$\psi'(\lambda) = \langle \nabla_h G(h(\lambda)), h'(\lambda) \rangle \quad (24)$$

At $\lambda = 0$, use $h'(0)$ from Eq. 19:

$$\psi'(0) = \langle \nabla_h G(h_S^{\text{KD}}), -(H_{J_{\text{KD}}}(h_S^{\text{KD}}))^{-1} \nabla_h G(h_S^{\text{KD}}) \rangle \quad (25)$$

Let $v = \nabla_h G(h_S^{\text{KD}})$. Then: $\psi'(0) = -\langle v, (H_{J_{\text{KD}}})^{-1} v \rangle$.

Since J_{KD} is μ -strongly convex, $H_{J_{\text{KD}}}$ is positive definite. Its inverse is also positive definite.

So if $v \neq 0$, then $\langle v, (H_{J_{\text{KD}}})^{-1} v \rangle > 0$. Thus, $\psi'(0) < 0$.

By continuity and the Mean Value Theorem (as in Lemma 1), for small $\lambda > 0$: $G(h(\lambda)) < G(h(0))$. □

Theorem 3 (Comparative Generalization Bound). *Under the stated assumptions, with probability at least $1 - \delta$:*

$$\begin{aligned} R(h_S^{\text{CCKD}}) &\leq R(h_S^{\text{KD}}) - \lambda \Delta_G - \frac{\mu}{2} \Delta_{\|\cdot\|^2} \\ &\quad + \frac{C}{\sqrt{n}} (\|h_S^{\text{CCKD}}\| + \|h_S^{\text{KD}}\|) + 4L_R \sqrt{\frac{\ln(2/\delta)}{2n}} \end{aligned} \quad (26)$$

Here: $-\Delta_G = G(h_S^{\text{KD}}) - G(h_S^{\text{CCKD}})$, $-\Delta_{\|\cdot\|^2} = \|h_S^{\text{KD}}\|^2 - \|h_S^{\text{CCKD}}\|^2$, $C = 2L_R \sqrt{K \cdot \max_i \Theta_S(\mathbf{x}_i, \mathbf{x}_i)}$ (independent of n).

Detailed Proof. Step 1: Rademacher Complexity Bounds.

Let $\mathcal{F}_B = \{h \in \mathcal{H}_{\Theta} : \|h\|_{\mathcal{H}_{\Theta}} \leq B\}$.

With probability at least $1 - \delta/2$, for any $h \in \mathcal{F}_B$:

$$R(h) \leq \hat{R}_n(h) + 2L_R \hat{\mathfrak{R}}_n(\mathcal{F}_B) + \text{const} \cdot \sqrt{\frac{\ln(2/\delta)}{n}}$$

In RKHS, $\hat{\mathfrak{R}}_n(\mathcal{F}_B) \leq \frac{B\sqrt{\text{tr}(\Theta_S)}}{n}$ [4].
Typically, $\text{tr}(\Theta_S) \leq n \cdot \max_i \Theta_S(\mathbf{x}_i, \mathbf{x}_i)$.
So, the complexity term is:

$$2L_R \frac{B\sqrt{n \cdot \max_i \Theta_S(\mathbf{x}_i, \mathbf{x}_i)}}{n} = \frac{2L_R B \sqrt{\max_i \Theta_S(\mathbf{x}_i, \mathbf{x}_i)}}{\sqrt{n}} \quad (27)$$

Define $C = 2L_R \sqrt{K \cdot \max_i \Theta_S(\mathbf{x}_i, \mathbf{x}_i)}$. (The \sqrt{K} comes from multi-class extension.)

Replace fixed B with actual solution norms. By union bound, with probability $\geq 1 - \delta$:

$$R(h_S^{\text{CCKD}}) \leq \hat{R}_n(h_S^{\text{CCKD}}) + \frac{C}{\sqrt{n}} \|h_S^{\text{CCKD}}\| + 2L_R \sqrt{\frac{\ln(2/\delta)}{2n}} \quad (28)$$

$$\hat{R}_n(h_S^{\text{KD}}) \leq R(h_S^{\text{KD}}) + \frac{C}{\sqrt{n}} \|h_S^{\text{KD}}\| + 2L_R \sqrt{\frac{\ln(2/\delta)}{2n}} \quad (29)$$

Step 2: CCKD Optimality Condition.

By optimality of h_S^{CCKD} :

$$\hat{R}_n(h_S^{\text{CCKD}}) \leq \hat{R}_n(h_S^{\text{KD}}) - \lambda \Delta_G - \frac{\mu}{2} \Delta_{\|\cdot\|^2} \quad (30)$$

Step 3: Combining the Inequalities.

Start with Eq. 28. Plug in Eq. 30. Then plug in Eq. 29. Result:

$$\begin{aligned} R(h_S^{\text{CCKD}}) &\leq \hat{R}_n(h_S^{\text{KD}}) - \lambda \Delta_G - \frac{\mu}{2} \Delta_{\|\cdot\|^2} \\ &\quad + \frac{C}{\sqrt{n}} \|h_S^{\text{CCKD}}\| + 2L_R \sqrt{\frac{\ln(2/\delta)}{2n}} \\ &\leq R(h_S^{\text{KD}}) + \frac{C}{\sqrt{n}} \|h_S^{\text{KD}}\| + 2L_R \sqrt{\frac{\ln(2/\delta)}{2n}} \\ &\quad - \lambda \Delta_G - \frac{\mu}{2} \Delta_{\|\cdot\|^2} + \frac{C}{\sqrt{n}} \|h_S^{\text{CCKD}}\| \\ &\quad + 2L_R \sqrt{\frac{\ln(2/\delta)}{2n}} \end{aligned} \quad (31)$$

Collect terms:

$$\begin{aligned} R(h_S^{\text{CCKD}}) &\leq R(h_S^{\text{KD}}) - \lambda \Delta_G - \frac{\mu}{2} \Delta_{\|\cdot\|^2} \\ &\quad + \frac{C}{\sqrt{n}} (\|h_S^{\text{KD}}\| + \|h_S^{\text{CCKD}}\|) + 4L_R \sqrt{\frac{\ln(2/\delta)}{2n}} \end{aligned} \quad (32)$$

□

Theorem 4 (Sufficient Condition for Superior Generalization Bound). *For large enough n , there is a non-empty range of λ . Choose λ from this range. Then, the generalization bound for h_S^{CCKD} is strictly tighter than for h_S^{KD} .*

Detailed Proof. From Theorem 3, the bound for h_S^{CCKD} is tighter if:

$$\lambda \Delta_G + \frac{\mu}{2} \Delta_{\|\cdot\|^2} > \frac{C}{\sqrt{n}} (\|h_S^{\text{CCKD}}\| + \|h_S^{\text{KD}}\|) + 4L_R \sqrt{\frac{\ln(2/\delta)}{2n}} \quad (33)$$

We compare the **order** of both sides in λ and n .

Left-Hand Side (LHS) Analysis:

We study $\lambda \Delta_G + \frac{\mu}{2} \Delta_{\|\cdot\|^2}$ for small $\lambda > 0$. Use perturbation results from earlier lemmas.

- **Norm Reduction Term:** From Lemma 1, $\Delta_{\|\cdot\|^2} = \|h_S^{\text{KD}}\|^2 - \|h_S^{\text{CCKD}}\|^2 \approx -\lambda \phi'(0)$. Since $\phi'(0)$ is strictly negative (Assumption 3), $\Delta_{\|\cdot\|^2} = O(\lambda)$. So, $\frac{\mu}{2} \Delta_{\|\cdot\|^2} = O(\lambda)$.
- **Feature Alignment Term:** From Lemma 2, $\Delta_G = O(\lambda)$. Then $\lambda \Delta_G = \lambda \cdot O(\lambda) = O(\lambda^2)$.

Total LHS = $O(\lambda) + O(\lambda^2)$. For small λ , the $O(\lambda)$ term dominates. So, LHS = $O(\lambda)$.

Right-Hand Side (RHS) Analysis:

The norms $\|h_S^{\text{CCKD}}\|$, $\|h_S^{\text{KD}}\|$ are bounded. Constants C , L_R do not depend on n or λ .

So RHS = $O(n^{-1/2})$.

Condition for Improvement:

We need LHS \geq RHS:

$$O(\lambda) > O(n^{-1/2})$$

This means: $\lambda > c_1/\sqrt{n}$, for some $c_1 > 0$ (problem-dependent).

But our analysis only holds for $\lambda < \lambda_0$, for some fixed $\lambda_0 > 0$.

For large n , c_1/\sqrt{n} becomes smaller than λ_0 .

So, the interval $(\frac{c_1}{\sqrt{n}}, \lambda_0)$ is non-empty.

Pick any λ in this interval.

Then, the CCKD improvement outweighs the statistical error. The bound becomes strictly tighter. □

A.4 Discussion and Interpretation

Our analysis gives a formal theory for the generalization gains of CCKD under the NTK framework. But its scope is limited by assumptions and bound tightness. We now discuss key interpretations and limits.

1. **On Assumption 3.** This assumption is central to our main result on norm reduction (Lemma 1), yet its validity is non-trivial. It requires a positive correlation between h_S^{KD} and the pre-conditioned gradient $(H_{J_{\text{KD}}})^{-1} \nabla_h G(h_S^{\text{KD}})$. The operator $(H_{J_{\text{KD}}})^{-1}$ effectively amplifies components of the gradient corresponding to flat directions in the KD loss landscape. Thus, the assumption implies that CCKD is most effective when the feature alignment term aligns with low-curvature directions of the knowledge distillation

objective. Empirically verifying this geometric condition in high-dimensional, practical models remains an open challenge.

2. **The λ versus n trade-off.** Theorem 4 reveals an important trade-off: the choice of λ cannot be made independently of the sample size n . Specifically:

- The perturbation analysis in Lemmas 1 and 2 is only valid for small λ , requiring $\lambda < \lambda_0$;
- To achieve a provable improvement in the generalization bound, the $O(\lambda)$ improvement from CCKD must dominate the $O(n^{-1/2})$ estimation error, implying $\lambda > c_1/\sqrt{n}$.

For sufficiently large n , the interval $(c_1/\sqrt{n}, \lambda_0)$ is non-empty, thereby defining a "sweet spot" for λ that yields a tighter generalization bound. This also explains the sensitivity of λ observed in practice:

- If λ is too small, the feature alignment signal is obscured by noise;
- If λ is too large, the solution deviates significantly from the locally convex region around h_S^{KD} , violating the assumptions of our theoretical framework.

3. **Limitation: Applicability of the NTK regime.** All our results rely on Assumption 1, which requires the neural tangent kernel Θ to remain fixed during training. This setting is representative of wide networks trained with small learning rates. However, modern deep networks often operate in a feature learning regime, where the kernel evolves substantially throughout training. Extending our analysis to account for such time-varying kernels remains a significant open problem. Currently, our theoretical guarantees are most directly applicable to overparameterized models operating in the NTK-like regime.

4. **Tightness of the generalization bounds.** Our analysis relies on Rademacher complexity bounds, which are known to be loose in practice and often fail to capture the precise generalization error of deep models. Therefore, the value of Theorems 3 and 4 is primarily qualitative rather than quantitative. They serve to elucidate a mechanistic explanation for the effectiveness of CCKD: by reducing the RKHS norm and enhancing feature alignment, CCKD provably leads to better generalization than standard KD within our theoretical framework.

References

- [1] R. Anil, G. Pereyra, A. Passos, R. Ormandi, G. E. Dahl, and G. E. Hinton. Large scale distributed neural network training through online distillation. *arXiv preprint arXiv:1804.03235*, 2018. 2, 3
- [2] N. Aronszajn. Theory of reproducing kernels. *Transactions of the American mathematical society*, 68(3):337–404, 1950. 9
- [3] J. Ba and R. Caruana. Do deep nets really need to be deep? In *Advances in neural information processing systems*, pages 2654–2662, 2014. 2
- [4] P. L. Bartlett and S. Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of machine learning research*, 3(Nov):463–482, 2002. 10, 11
- [5] M. Belkin, D. Hsu, S. Ma, and S. Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019. 9
- [6] S. P. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge university press, 2004. 10
- [7] C. Buciluă, R. Caruana, and A. Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '06, pages 535–541, New York, NY, USA, 2006. ACM. 3
- [8] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. Ieee, 2009. 6
- [9] J. Deng, J. Guo, and S. Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. 2018. 5, 7
- [10] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, and R. Fergus. Exploiting linear structure within convolutional networks for efficient evaluation. In *Advances in neural information processing systems*, pages 1269–1277, 2014. 1, 2
- [11] A. L. Dontchev and R. T. Rockafellar. *Implicit functions and solution mappings*, volume 543. Springer, 2009. 9
- [12] S. S. Du, X. Zhai, B. Poczos, and A. Singh. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018. 3
- [13] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *The Annals of statistics*, 32(2):407–499, 2004. 9
- [14] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014. 1
- [15] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. pages 87–102, 2016. 7
- [16] S. Han, H. Mao, and W. J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and Huffman coding. *arXiv preprint arXiv:1510.00149*, 2015. 2
- [17] S. Han, J. Pool, J. Tran, and W. Dally. Learning both weights and connections for efficient neural network. In *Advances in neural information processing systems*, pages 1135–1143, 2015. 1, 2
- [18] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5, 6

- [19] B. Heo, M. Lee, S. Yun, and J. Y. Choi. Improving knowledge distillation with supporting adversarial samples. *arXiv preprint arXiv:1805.05532*, 2018. [3](#)
- [20] B. Heo, M. Lee, S. Yun, and J. Y. Choi. Knowledge distillation with adversarial samples supporting decisionboundary. *arXiv preprint arXiv:1805.05532*, 2018. [3](#)
- [21] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. [1](#), [2](#), [3](#), [5](#), [6](#), [9](#)
- [22] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. [1](#), [2](#)
- [23] I. Hubara, M. Courbariaux, D. Soudry, E. Y. Ran, and Y. Bengio. Quantized neural networks: Training neural networks with low precision weights and activations. *Journal of Machine Learning Research*, 18, 2016. [2](#)
- [24] A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018. [9](#)
- [25] M. Jaderberg, A. Vedaldi, and A. Zisserman. Speeding up convolutional neural networks with low rank expansions. *arXiv preprint arXiv:1405.3866*, 2014. [1](#), [2](#)
- [26] I. Kemelmachershizerman, S. M. Seitz, D. Miller, and E. Brossard. The megaface benchmark: 1 million faces for recognition at scale. In *Computer Vision and Pattern Recognition*, pages 4873–4882, 2016. [7](#)
- [27] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009. [5](#)
- [28] J. Lee, L. Xiao, S. Schoenholz, Y. Bahri, R. Novak, J. Sohl-Dickstein, and J. Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. *Advances in neural information processing systems*, 32, 2019. [9](#)
- [29] Q. Li, S. Jin, and J. Yan. Mimicking very efficient network for object detection. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7341–7349. IEEE, 2017. [1](#), [3](#), [5](#), [7](#)
- [30] T.-Y. Lin, A. RoyChowdhury, and S. Maji. Bilinear cnn models for fine-grained visual recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1449–1457, 2015. [4](#)
- [31] P. Luo, Z. Zhu, Z. Liu, X. Wang, X. Tang, et al. Face model compression by distilling knowledge from neurons. In *AAAI*, pages 3560–3566, 2016. [5](#), [7](#)
- [32] P. Molchanov, S. Tyree, T. Karras, T. Aila, and J. Kautz. Pruning convolutional neural networks for resource efficient inference. 2016. [1](#), [2](#)
- [33] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014. [2](#), [3](#)
- [34] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4510–4520, 2018. [5](#)
- [35] B. B. Sau and V. N. Balasubramanian. Deep model compression: Distilling knowledge from noisy teachers. *arXiv preprint arXiv:1610.09650*, 2016. [2](#)
- [36] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [1](#)
- [37] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015. [1](#)
- [38] G. Urban, K. J. Geras, S. E. Kahou, O. Aslan, S. Wang, R. Caruana, A. Mohamed, M. Philipose, and M. Richardson. Do deep convolutional nets really need to be deep and convolutional? *Nature*, 521, 2016. [3](#)
- [39] G. Ver Steeg and A. Galstyan. Discovering structure in high-dimensional data through correlation explanation. In *Advances in Neural Information Processing Systems*, pages 577–585, 2014. [4](#)
- [40] F. Wang, L. Chen, C. Li, S. Huang, Y. Chen, C. Qian, and C. C. Loy. The devil of face recognition is in the noise. *arXiv preprint arXiv:1807.11649*, 2018. [7](#)
- [41] L. Wei, S. Zhang, W. Gao, and Q. Tian. Person transfer gan to bridge domain gap for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 79–88, 2018. [6](#)
- [42] J. Wu, L. Cong, Y. Wang, Q. Hu, and J. Cheng. Quantized convolutional neural networks for mobile devices. In *Computer Vision and Pattern Recognition*, pages 4820–4828, 2016. [2](#)
- [43] Z. Xu, Y.-C. Hsu, and J. Huang. Training shallow and thin networks for acceleration via knowledge distillation with conditional adversarial networks. 2018. [3](#), [5](#), [7](#)
- [44] J. Yim, D. Joo, J. Bae, and J. Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, 2017. [2](#), [3](#)
- [45] S. Zagoruyko and N. Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. 2016. [2](#), [3](#), [5](#), [6](#)
- [46] X. Zhang, X. Zhou, M. Lin, and J. Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. 2017. [1](#), [2](#)
- [47] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu. Deep mutual learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4320–4328, 2018. [2](#), [5](#), [6](#)