

Knowledge Distillation as Efficient Pre-training: Faster Convergence, Higher Data-efficiency, and Better Transferability

Ruifei He^{1*}, Shuyang Sun^{2†}, Jihan Yang^{1†}, Song Bai^{3‡}, Xiaojuan Qi^{1‡}

¹The University of Hong Kong ²University of Oxford ³ByteDance

{ruifeihe, jhyang, xjqqi}@eee.hku.hk, kevinsun@robots.ox.ac.uk, songbai.site@gmail.com

Abstract

Large-scale pre-training has been proven to be crucial for various computer vision tasks. However, with the increase of pre-training data amount, model architecture amount, and the private/inaccessible data, it is not very efficient or possible to pre-train all the model architectures on large-scale datasets. In this work, we investigate an alternative strategy for pre-training, namely Knowledge Distillation as Efficient Pre-training (**KDEP**), aiming to efficiently transfer the learned feature representation from existing pre-trained models to new student models for future downstream tasks. We observe that existing Knowledge Distillation (KD) methods are unsuitable towards pre-training since they normally distill the logits that are going to be discarded when transferred to downstream tasks. To resolve this problem, we propose a feature-based KD method with non-parametric feature dimension aligning. Notably, our method performs comparably with supervised pre-training counterparts in 3 downstream tasks and 9 downstream datasets requiring **10×** less data and **5×** less pre-training time. Code is available at <https://github.com/CVMI-Lab/KDEP>.

1. Introduction

With the booming of large-scale datasets [16, 46, 37, 61, 57], many computer vision tasks have benefitted significantly from pre-training in the past decade. In fact, it has been a de facto strategy to first pre-train on datasets like ImageNet [16] and then fine-tune on downstream tasks [75, 8, 24, 68, 55], especially when the data of downstream tasks is scarce.

However, the increasing pre-training data scale and the inaccessibility of private data [61] render pre-training all architectures on large datasets not efficient or possible. As well-trained deep neural networks are essentially condensed

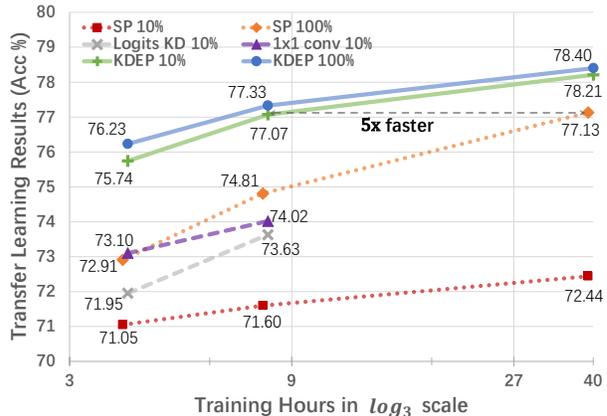


Figure 1. Transfer performance (averaged top-1 accuracy of four image classification tasks (details in Sec. 4)) compared to Supervised Pre-training (SP), traditional KD method (logits KD, 1×1 conv), and KDEP (SVD+PTS) with different data amount (10% or 100% ImageNet-1K data) and training schedules.

memory bank of datasets [3, 20], we wonder whether the condensed data knowledge encoded into a pre-trained model can be leveraged to efficiently pre-train new architectures with only a relatively small set of pre-training data?

In this work, we propose Knowledge Distillation as Efficient Pre-training (KDEP), transferring the feature extraction capability of the teacher obtained from large-scale data, to the student model for solving future downstream tasks. Note that KDEP is quite different from traditional Knowledge Distillation (KD) that only targets at distilling the knowledge of a given specific task to a student model.

Studies of existing KD methods for KDEP. Our empirical studies show that existing KD methods such as logits KD [29] (*i.e.* distilling the task-specific output logits) and feature-level KD [27] lead to inferior performance (see Figure 1: “logits KD” and “ 1×1 conv”), indicating that existing KD methods tailored to different tasks might be unable to fully leverage the knowledge condensed in the teacher model when pre-training new models with limited data and computation budget.

After further investigation, we conclude a potential issue

*Part of the work is done during an internship at ByteDance AI Lab.

†Equal contribution.

‡Corresponding author.

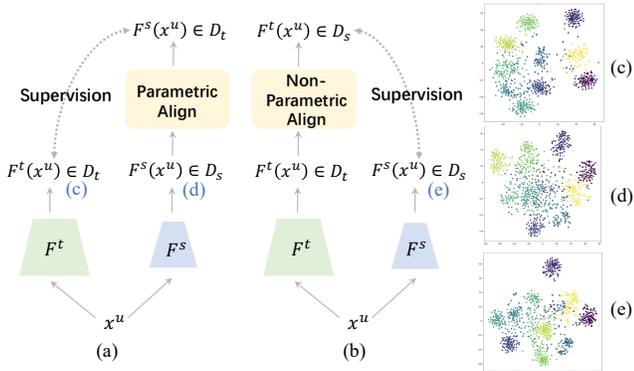


Figure 2. **Framework and visualization.** Here we illustrate the KDEP framework of feature-based KD with (a): parametric aligning; and (b): non-parametric aligning. Notations refer to Sec. 3.1. Notice that only the learned F^s is to transfer to downstream tasks. Moreover, we visualize the feature representation of (c): original teacher, (d): student distilled by (a) with 1×1 conv aligning, and (e): student distilled by (b) with our SVD+PTS aligning. For visualization, we randomly sample 10 classes in ImageNet-1K and use 100 samples from each class. T-SNE [64] is used for dimensionality reduction.

of “indirect feature learning”, where the supervision of distillation is not *directly applied* on the feature extractor that will be transferred, but on a new learnable module added after it and jointly optimized with it, which turns the feature representation learning into an indirect process. Specifically, “logits KD” applies the supervision on the classifier’s output, where the learnable parameters of the classifier is jointly learned with the feature extractor. Although feature-level KD applies supervisions on the features, most feature-level KD methods adopt a parametric module to align the feature dimensions of teacher and student, typically a 1×1 convolutional layer (“ 1×1 conv”) [58, 27, 28, 2], which again adds learnable parameters after the feature extractor and forms an indirect process. As shown in Figure 2d, we visualize the learned feature representation of feature-based KD with 1×1 conv alignment: the learned feature representation fails to follow that of the teacher’s (Figure 2c). Consequently, the pre-trained models only deliver suboptimal transfer learning performance.

Our Method. Motivated by the identified potential problem, our KDEP investigates non-parametric methods for aligning the feature dimensions to avoid indirect feature learning. Empirically, we found Singular Value Decomposition (SVD) works effectively by compressing¹ the features with minimal information loss.

However, features processed by SVD will trigger the component domination effect [13], *i.e.* the feature variances among channels are of great magnitude differences, and

¹We focus on the setting that teachers having larger dimensions in our study, since this is more frequently met in real-world applications. Foundation models tend to have larger feature dimensions where models for deployment usually have smaller dimensions.

largely differ from those of normal DCNNs. This interferes the optimization of the network. To further boost feature learning, we design a Power Temperature Scaling (PTS) function to reduce the variance differences while preserving the original relative magnitude, which tailors features from SVD for DCNNs. As illustrated in Figure 2e, with our SVD+PTS aligning method, the distilled student obtains feature representation similar to the original teacher’s (Figure 2c) while matching the feature dimensions. Notably, our method adds no learnable parameters, does not rely on the task loss or the logits loss [29], and only use supervision from the teacher’s penultimate feature (after global average pooling), which is more general for feature representation learning and allows more potential pre-trained teachers.

Results. Our major findings of KDEP are summarized in Figure 1: 1) **Faster convergence.** Our method achieves comparable or better transfer learning results with only 10% or 20% training time than supervised pre-training (SP) on the whole ImageNet-1K dataset. 2) **Higher data-efficiency.** With only 10% of ImageNet-1K unlabeled data (discard the labels) and an available pre-trained teacher model, our distilled student obtains better transfer learning results than SP with 100% ImageNet-1K data. 3) **Better Transferability.** Given the same computation budget and data amount as SP, our method achieves higher transfer learning performance. With the proposed KDEP method, we could realize pre-training once and distilling it to all: distilling a pre-trained teacher (either to utilize an available pre-trained model or pre-train a certain architecture) to efficiently pre-train all other student models.

2. Related Works

Transfer Learning (TL), usually by fine-tuning a pre-trained model to a downstream task with labeled data, has become a common practice in machine learning problems and applications. To better understand TL, it can be separated into two steps: pre-training and transfer.

Recent years have witnessed increasing successful works on pre-training, including supervised pre-training [32, 39, 46, 61, 34], self-supervised pre-training [9, 23, 5, 22, 11, 74, 70], and semi-supervised pre-training [66, 54, 10]. Given the pre-trained models, the next step is to transfer the learned representation to a target task. Except from the widely used fine-tuning method [72, 1], there are also other methods proposed for better exploiting the knowledge absorbed in pre-training, such as L2-SP [67], DELTA [40], BSS [12], and Co-Tuning [73].

While large-scale pre-training yield better representation and downstream performance, the cost of pre-training is also rapidly increasing [61, 17, 34]. Therefore, we hope to propose an orthogonal strategy for pre-training, distilling a pre-trained model to pre-train different student models.

Knowledge Distillation (KD) has been developed as an ef-

fective way for model compression and acceleration, and various methods mainly falls into three research streams: response-based [4, 29, 77], feature-based [28, 58, 35, 30, 50, 33, 31, 6, 27], and relation-based [71, 38, 51, 44, 49, 7, 63, 53, 62] methods.

Response-based methods usually take the final outputs called logits as supervision and use a temperature factor to adjust the smoothness [29]. Though the logits introduce “dark knowledge” into training and show improved results, response-based KD leaves out large information of the intermediate features, which are found to be crucial for representation learning [58].

Feature-based KD was first introduced by Fitnets [58], using intermediate feature maps as hints to improve KD performance. Following Fitnets, attention maps [35], neuron selectivity patterns [30], paraphraser [33], route constraint [31], and activation boundary [28] are also proposed to better utilize feature-level knowledge. Heo *et al.* [27] investigate different design aspects of feature-based KD, and propose margin ReLU, Pre-ReLU distillation position, and a partial L2 loss function. Chen *et al.* [6] use an attention mechanism to adaptively assign proper teacher layers for each student layer.

Relation-based methods further boost the performance by utilizing the relationships between different layers or data samples. FSP [71] uses inner products between features of two layers as a flow of solution process, but it restricts to the same feature map sizes between teacher and student. Lee *et al.* [38] also utilize the correlations between feature maps by using Radial Basis Function, and apply SVD on spatial dimensions to both teacher and student’s feature maps to avoid the mismatch in spatial resolutions. However, they still need the dimension of feature maps to be equal for teacher and student. In contrast, our method can adapt to teacher and student pairs with different feature map resolutions and dimensions. Another line in relation-based methods is utilizing the relation between data samples, where different mechanisms have been proposed, including instance relationship graph [44], similarity matrix [63], similarity probability distribution [51] and so on.

However, traditional KD methods only focus on a single task and transfer task-specific knowledge, while our KDEP focuses on the transferability of the distilled student, which distinguishes us from most previous KD methods. To our best knowledge, only Li *et al.* [41] try to utilize KD for student’s transferability. They show traditional KD would hurt the transferability of the student, and propose a multi-head, multi-task distillation method using an unlabeled proxy dataset and a generalist teacher to improve downstream performance of the distilled student.

Nevertheless, their method needs the student model initialized by ImageNet pre-trained weights, multiple teachers fine-tuned on the domains which are related to the down-

stream domains, and a multi-head, multi-task training procedure, which violates our efficient pre-training setting. On the contrary, our method optimizes the student from scratch, and only needs a single generalist teacher, an unlabeled dataset with a simple yet efficient training pipeline to achieve comparable transfer learning performance with supervised pre-training.

3. KDEP

3.1. Overview

We define the KDEP settings as below: given a teacher model F^t (pre-trained on a large-scale dataset \mathcal{D}) and a set of unlabeled examples $\mathcal{D}_u = \{x_i^u\}_{i=1}^{N_u}$ (N_u is the number of unlabeled images), our goal is to pre-train a student model F^s to generalize well on various downstream tasks. Note that the dataset scale of \mathcal{D}_u could be magnitude smaller than \mathcal{D} and only \mathcal{D}_u is available during the student training. Since we focus on feature representation learning rather than tailoring the model to a specific task, both F^t and F^s yield the feature representations instead of task-specific logits. We denote the shape of $F^t(x_i^u)$, $F^s(x_i^u)$ as D_t, D_s . The training objective of the KDEP method is

$$\frac{1}{N_u} \sum_{i=1}^{N_u} \mathcal{L}(F^t(x_i^u), F^s(x_i^u)), \quad (1)$$

where \mathcal{L} is the L_2 loss. To meet the demand of our proposed KDEP, several under-explored obstacles are needed to be addressed.

The first is a known issue for feature-based KD, that is the feature dimension mismatch (*i.e.*, $D_t \neq D_s$) between the teacher and student. In our study, we found the frequently adopted strategy, to add a parametric module like a 1×1 conv, is sub-optimal for our KDEP settings. Instead, we demonstrate that non-parametric methods (*e.g.* SVD) are more effective than 1×1 conv for aligning the dimensions. Analysis and details will unfold in Sec.3.2.

The second is a byproduct of our non-parametric feature dimension aligning method: the feature statistics after the alignment would differ from normal DCNNs. Hence, we study several mechanisms for correcting the feature statistics and conclude them as a transformation module. We will further elaborate on the design choices in Sec. 3.3.

The third is still an open issue even after our exploration: What is a good teacher for KDEP? Our empirical studies show that stronger models are not necessarily better teachers, and we find the compactness of the teacher’s feature distribution to be a crucial indicator (see Sec. 3.4). We hope we could inspire more future works on this topic.

3.2. Aligning Feature Dimensions

Motivated by the indirect feature learning problem, we propose non-parametric feature dimension aligning methods with several variants. Concretely, previous parametric

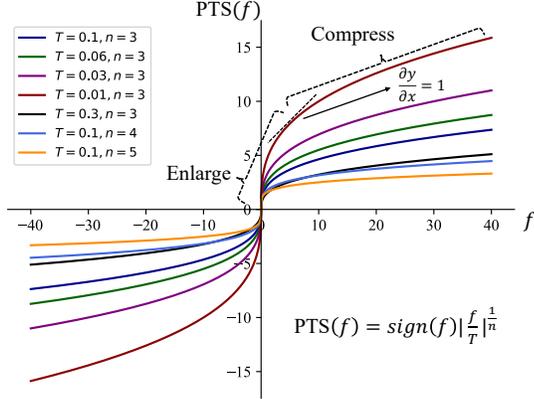


Figure 3. We propose PTS function for reducing Std Ratio while preserving the original relative magnitude. Curves of different values of T and n are shown.

methods add a parametric module to align D_s toward the supervision D_t (ref. Figure 2a). In contrast, as presented in Figure 2b, we apply a non-parametric method to project D_t to D_s , which can then serve as the supervision directly.

We investigate three variants for non-parametric aligning methods: channel selection, interpolation, and SVD. Along this line, SVD stands out thanks to its power of effectively compressing the feature-level knowledge and maintaining as much information as possible. Detailed experimental results would be included in Sec. 4.

Pre-ReLU distillation feature position has been used in previous feature-based KD [28, 27] and shown improved results. In our methods, we distill the features before the ReLU activation function for one more consideration, that SVD’s outputs contain both negative and positive values.

3.3. Transformation Module

While SVD effectively compresses the features with unnoticeable information loss, it brings along difficulties for optimization. After the SVD alignment, the feature variances of different channels are of magnitude differences, whereas those of normal DCNNs are usually within the same magnitude. Concretely, we define the term Std Ratio as the largest standard deviation (Std) to the smallest Std among all feature channels at the penultimate features across all training data samples. According to our study, we found the Std Ratio of features after SVD is usually over $10\times$ larger than that of normal DCNNs.

As a result, the L_2 loss tends to be dominated by the feature channels with the largest variances, leaving the minor ones under-fitted, for which we provide theoretical analysis below. We view the value of each feature channel from the teacher as a random variable (T), which is of zero mean after SVD. Also, since the student is optimized from random initialization [25], we regard each feature channel of the student also as a random variable (S) with zero mean. As shown in Theorem 1 and Proof 1, we state and prove

that the mathematical expectation of the L_2 loss from each feature channel increases monotonically with the Std of the teacher’s feature channel, which explains the difficulty of learning from teacher with large Std Ratio.

Theorem 1 Given two independent random variables with normal distribution $T \sim N(0, \sigma^2)$ and $S \sim N(0, \sigma_s^2)$, then $F(\sigma) = \mathbb{E}[(T - S)^2]$ is monotonically increasing ($\sigma > 0$).

Proof 1 (Detailed Proof in Appendix.)

$$\begin{aligned}
 F(\sigma) &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} (t - s)^2 \cdot P(t, s) dt ds \\
 &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} e^{-\frac{t^2}{2\sigma^2}} (t^2 + \sigma_s^2) dt = \sigma^2 + \sigma_s^2 \\
 \frac{dF(\sigma)}{d\sigma} &= 2\sigma > 0 \Rightarrow \text{monotonically increasing}
 \end{aligned}$$

Therefore, we propose to reduce the Std Ratio after SVD to normal ranges of DCNNs by a transformation module. A simple method of Scale Normalization (SN) has been used in previous works [13, 38], which divides each feature channel by its corresponding Std to ensure each channel has similar scales. Similarly, we also experiment with a variant of SN: rather than dividing the corresponding Std to obtain similar scales, we scale each channel’s Std to match the top- D_s Std of features before SVD (target Std), which we name as Std Matching (SM). However, both SN and SM are local transformations that transform channel-wisely, and thus may fail to preserve the original relative magnitude between different feature channels (example in Appendix).

To match statistics without hurting the original relative magnitude, we propose to use a global non-decreasing transformation function that can reduce the Std Ratio while maintaining the relative magnitude. Concretely, we control the value ranges by a temperature parameter T similar to logits KD [29], and then apply a power operation while preserving the signs. We refer to the function as Power Temperature Scaling (PTS), which is as follows:

$$\text{PTS}(f) = \text{sign}(f) \left| \frac{f}{T} \right|^{\frac{1}{n}}, \quad (2)$$

where n is the parameter of the exponent and f is the input value. As shown in Figure 3, the PTS function could enlarge small values and compress large values, while globally non-decreasing, and thus fulfills our goal of matching normal statistics and preserving the relative magnitude.

3.4. Teacher Selection

In this section, we further explore how to select a good teacher for KDEP. Naturally, we would consider utilizing stronger models as the teacher, where we study and compare several paradigms of potential stronger models:

- **Standard SP**: the most frequent supervised pre-training strategy that pre-trains a model on ImageNet-1K; we use the architecture of ResNet50 (R50) [26] as teacher for this study;

- SP with more data: we experiment with an available Microsoft (MS) Vision R50² pre-trained with four datasets (over 40 million data): ImageNet-22K, COCO, and two webly-supervised datasets [69]. The provided weights only contain the feature extractor;
- Pre-training with unlabeled data: we explore a semi-weakly supervised pre-trained R50 (SWSL);
- Distilled models: we use a MEAL V2 R50 trained by distillation on ImageNet-1K.
- Advanced architecture: we experiment with a state-of-the-art architecture Swin Transformer [45] and use Swin-B pre-trained with ImageNet-22K and fine-tuned on ImageNet-1K.

We empirically found that stronger models (*i.e.* higher performance on ImageNet-1K benchmark) do not necessarily achieve better KDEP performance (*i.e.* distilled student’s transferability under the KDEP setting), which resonates with previous findings in KD [48, 47] that more accurate teachers may distill worse. To investigate the reasons, we visualize the feature representation of different teachers (ref. Figure 4), and surprisingly found that the KDEP performance has a strong correlation with the teacher’s feature compactness (compactness means the feature distribution of data samples of the same class lying tightly in the feature space while those of different classes far from each other). Detailed analysis and results are in Sec. 4.3.

4. Experiments

4.1. Experimental Setup

Datasets and downstream tasks. For the proposed KDEP settings, we use the ImageNet-1K [16] dataset as unlabeled data by abandoning the labels, and we use 10% or 100% of the dataset for different settings.

To evaluate the transfer learning performance of the models, we evaluate on three popular downstream tasks: image classification, semantic segmentation, and object detection. For image classification, we select four diverse datasets to study the transferability: CIFAR-100 [36], CUB-200-2011 [65], DTD [14], and Caltech-256 [21]. For semantic segmentation, we use three widely used datasets: Cityscapes [15], PASCAL VOC 2012 (VOC12) [18], and ADE20K [76]. For object detection, we evaluate the transfer performance on two benchmarks: PASCAL VOC [19], and COCO [43].

Teacher-Student (T-S) pair. We experiment with two different teacher-student pairs, R50 \rightarrow ResNet18 (R18), and R50 \rightarrow MobileNetV2 (MNv2) [60], representing knowledge transfer between similar and dissimilar networks, respectively. For the teacher model, we conduct our main experiments with standard SP R50 and MS R50.

Comparison methods. We mainly compare KDEP with supervised pre-training (SP). We denote the student supervised pre-trained with all ImageNet-1K data for 90 epochs as SP oracle (SP. o.), and with fewer data or for shorter schedule as SP baseline (SP. b.) in each setting.

Implementation details. We implement our method using PyTorch [52] and all experiments are conducted using four 32G V100 GPUs. For studying KDEP, we explore different settings with various data amounts and training schedules. For the 10% ImageNet-1K data setting, we set the training epochs to 90 or 180; when using 100% ImageNet-1K data, we train for 9 or 18 epochs to verify fast convergence, and for 90 epochs to further boost performance, where 90 epochs with all ImageNet-1K data is the standard supervised pre-training schedule [59]. For all downstream tasks, we use the same schedule and evaluation protocols for all models for a fair comparison. More elaborated implementation details are given in the Appendix.

Evaluation. We report the top-1 accuracy, mean Intersection over Union (mIoU), and AP, AP₅₀, AP₇₅ for classification, segmentation, and detection, respectively. All results are averaged over at least three trials. Time refers to the pre-training time of SP or KDEP on four 32G V100 GPUs.

4.2. Main Results

In this section, we compare our best KDEP method (SVD+PTS) with supervised pre-training under different data amounts and training schedules. We evaluate the transferability on all 9 transfer learning tasks covering image classification, semantic segmentation and object detection. With our extensive experimental results (*e.g.* Table 1), we demonstrate that Knowledge Distillation can be used as an effective way of pre-training, outperforming standard supervised pre-training with fewer training data and shorter training schedules. In the following, we explore the transferability, data-efficiency, and convergence speed of KDEP respectively under different setups.

Note that all KDEP methods use the MS R50 as the teacher in this section. Results with more teachers are in our ablation studies. Also, due to the length limit, we show the results of R18 as student in our paper, and MNv2 as student in the Appendix, where similar results are achieved.

Exploring transferability under 10% data with short schedules. In this setting, we use only 10% ImageNet-1K data, which is a total of 128k images randomly sampled from the original 1.28 million images. We pre-train for 90 epochs or 180 epochs with KDEP or SP. As shown in Table 1, KDEP significantly outperforms its SP counterparts in different settings.

We take the transfer performance in classification as an example to illustrate and analyze in this setting and the settings following unless noted. SP models’ performance drops significantly when pre-trained with only 10% data

²<https://pypi.org/project/microsoftvision/>

(77.13→71.05). Even increasing the training schedule to 180 epochs (71.05→71.60) or 900 epochs (71.05→72.44) fails to eliminate such performance drop. In contrast, with only 90 epochs, KDEP (75.74) largely bridges the gap between supervised pre-training baselines and oracle. Increasing the schedule to 180 epochs further closes the gap to an unnoticeable drop (77.07 vs. 77.13) while only using 10% data and around 20% training time. Similar results are also observed in segmentation and detection results as well as the R50 → MNV2 pair (in Appendix).

Exploring transferability under the same data amount and schedule as standard SP. We further explore the KDEP performance given the same data amount and training schedules as standard SP. In Table 1, with 100% data and 90 training epochs, KDEP produces models with much stronger transferability (78.40) than standard SP (77.13), while only adding minuscule computation costs.

Exploring data-efficiency. Here, we again only use 10% data but extend the training schedule to 900 epochs, which aligns the training iterations with the standard supervised pre-training setting. In Table 1, KDEP could consistently outperform the supervised oracle on all 9 tasks with only 10% ImageNet-1K data, which verifies our initial idea that models pre-trained on large-scale data could transfer the condensed data knowledge to other architectures even without using full pre-training data.

Exploring convergence speed. For further exploring convergence speed, we compare KDEP and SP under different data amounts (*i.e.* 10% and 100% ImageNet-1K data) as the training time increases. As shown in Figure 1, the transfer performance of SP increases almost uniformly with the training time, while KDEP shows a favourable characteristic of fast convergence under both data amounts. Notably, KDEP produces comparable or superior transfer performance as standard SP with 5× fewer training time.

4.3. Ablation Studies

For our KDEP settings, we notice a strong correlation between the transfer learning results of different tasks in Sec. 4.2. Hence, we use four image classification tasks to evaluate the transferability in our ablation studies. All results in our ablation studies are the averaged top-1 accuracy over four classification tasks.

Ablation study on feature dimension aligning methods. In this ablation study, we aim to investigate the effectiveness of various feature dimension aligning methods for KDEP. Firstly, we compare parametric methods with non-parametric methods under the 10% data and short schedule of 90 epochs setting with three Teacher-Student (T-S) pairs. For parametric methods, we experiment with three 1×1 conv variants: Post-ReLU, Pre-ReLU and the one in [27] (details in Appendix). For non-parametric methods, we explore channel selection (CS.var, CS.rand), interpola-

Method	Data	Epoch	Time (h)	Classification (Acc %)				
				Caltech	DTD	CUB	CIFAR	Avg
rand. init.	-	-	-	55.27	45.16	55.89	77.34	58.42
SP. b.	10%	90	3.9	68.83	66.17	69.93	79.29	71.05
KDEP	10%	90	4.0	75.33	71.80	74.20	81.61	75.74
SP. b.	100%	9	3.9	71.27	68.14	71.97	80.27	72.91
KDEP	100%	9	4.0	75.42	72.15	75.11	82.22	76.23
SP. b.	10%	180	7.8	70.09	66.14	70.84	79.34	71.60
KDEP	10%	180	8.0	77.15	72.67	75.99	82.23	77.07
SP. b.	100%	18	7.8	74.01	69.18	74.63	81.43	74.81
KDEP	100%	18	8.0	77.29	73.07	76.50	82.47	77.33
SP. b.	10%	900	39	71.10	67.02	72.16	79.49	72.44
KDEP	10%	900	40	79.00	74.28	76.89	82.64	78.21
SP. o.	100%	90	39	77.18	71.81	77.44	82.08	77.13
KDEP	100%	90	40	79.08	74.34	77.29	82.89	78.40

Method	Data	Epoch	Time (h)	Segmentation (mIoU %)			
				Cityscapes	VOC12	ADE20K	Avg
rand. init.	-	-	-	57.87	49.46	31.37	46.23
SP. b.	10%	90	3.9	69.56	68.30	34.39	57.41
KDEP	10%	90	4.0	70.41	72.34	36.09	59.61
SP. b.	100%	9	3.9	67.92	70.05	35.27	57.75
KDEP	100%	9	4.0	69.73	72.43	36.17	59.44
SP. b.	10%	180	7.8	69.89	69.79	35.03	58.24
KDEP	10%	180	8.0	70.27	72.82	36.60	59.90
SP. b.	100%	18	7.8	70.19	71.02	35.19	58.80
KDEP	100%	18	8.0	70.90	73.82	36.73	60.48
SP. b.	10%	900	39	69.85	69.55	35.52	58.31
KDEP	10%	900	40	71.93	74.28	37.30	61.17
SP. o.	100%	90	39	71.01	73.13	36.02	60.05
KDEP	100%	90	40	71.39	73.75	36.97	60.70

Method	Data	Epoch	Time (h)	Detection					
				VOC0712			COCO		
				AP	AP ₅₀	AP ₇₅	AP	AP ₅₀	AP ₇₅
rand. init.	-	-	-	26.7	52.5	23.1	25.2	41.9	26.3
SP. b.	10%	90	3.9	39.8	69.7	39.1	27.6	45.3	28.8
KDEP	10%	90	4.0	41.9	72.4	41.7	28.6	46.5	29.9
SP. b.	100%	9	3.9	40.4	70.5	40.2	27.5	45.3	28.6
KDEP	100%	9	4.0	42.5	73.3	43.3	28.8	46.9	30.1
SP. b.	10%	180	7.8	39.4	69.6	38.7	28.2	45.9	29.4
KDEP	10%	180	8.0	43.4	73.8	43.8	29.2	47.4	30.6
SP. b.	100%	18	7.8	41.2	71.5	40.7	28.1	46.1	29.3
KDEP	100%	18	8.0	43.3	73.6	44.2	29.3	47.5	30.8
SP. b.	10%	900	39	39.3	69.9	38.8	28.5	47.0	29.9
KDEP	10%	900	40	42.8	73.5	43.4	29.9	48.4	31.7
SP. o.	100%	90	39	41.8	72.6	41.6	29.0	47.3	30.4
KDEP	100%	90	40	42.8	73.9	43.4	29.7	48.2	31.3

Table 1. KDEP vs. SP, R50 → R18, fine-tuned on various tasks. KDEP refers to our SVD+PTS method. Note that COCO is used for the teacher’s pre-training while not for the student’s.

tion, and SVD (details in Appendix). As shown in Table 2, among three variants of 1×1 conv methods, the two that adopts Pre-ReLU feature distillation position give better performance. Interestingly, non-parametric methods consistently outperform all variants of the 1×1 conv method with significant gains. Moreover, SVD produces the best performance under various T-S pairs.

Further, we study the data-efficiency and convergence speed of parametric methods under various KDEP settings, with results shown in Table 3. For data-efficiency, parametric methods suffer even trained for 900 epochs with 10% data, showing low data-efficiency. For settings verifying

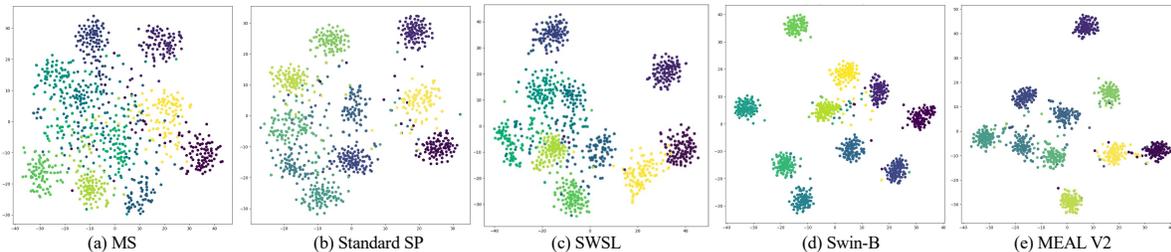


Figure 4. Visualization (same as in Figure 2) of feature distributions. Order ranked by compactness: diverse (left) to compact (right).

convergence speed (*i.e.* 100% data with 9 or 18 epochs), parametric methods perform similarly with the supervised pre-training baselines. The lacking of high data-efficiency and fast convergence characteristics may be caused by the added learnable module that hinders feature learning.

Teacher	Standard SP R50	MS R50	MS R50
Student	R18	R18	MNV2
1×1 conv (Post-ReLU)	71.37	72.43	71.12
1×1 conv (Pre-ReLU)	71.01	73.10	73.72
1×1 conv ([27])	71.35	72.98	73.62
CS.var	73.54	74.60	74.92
CS.rand	73.47	74.37	74.90
Interpolation	73.27	74.81	74.96
SVD	74.29	75.09	75.07

Table 2. **Ablation study on feature dimension aligning methods.** Setting: 10% data and 90 epochs. CS: “Channel Selection”.

Method	Data	Epoch	Time (h)	Acc (%)
SP. b.	10%	90	3.9	71.05
KDEP (1×1 conv, Pre-ReLU)	10%	90	4.0	73.10
KDEP (1×1 conv, [27])	10%	90	4.0	72.98
KDEP (SVD, PTS)	10%	90	4.0	75.74
SP. b.	100%	9	3.9	72.91
KDEP (1×1 conv, Pre-ReLU)	100%	9	4.0	73.24
KDEP (1×1 conv, [27])	100%	9	4.0	73.15
KDEP (SVD+PTS)	100%	9	4.0	76.23
SP. b.	10%	180	7.8	71.60
KDEP (1×1 conv, Pre-ReLU)	10%	180	8.0	74.02
KDEP (1×1 conv, [27])	10%	180	8.0	73.90
KDEP (SVD, PTS)	10%	180	8.0	77.07
SP. b.	100%	18	7.8	74.81
KDEP (1×1 conv, Pre-ReLU)	100%	18	8.0	74.32
KDEP (1×1 conv, [27])	100%	18	8.0	73.78
KDEP (SVD+PTS)	100%	18	8.0	77.33
SP. b.	10%	900	39	72.44
KDEP (1×1 conv, Pre-ReLU)	10%	900	40	75.29
KDEP (1×1 conv, [27])	10%	900	40	75.33
KDEP (SVD+PTS)	10%	900	40	78.21
SP. o.	100%	90	39	77.13

Table 3. **Ablation study on the 1×1 conv method.** T-S pair: MS R50 \rightarrow R18.

Ablation study on transformation module. We explore the effectiveness of three mechanisms of the transformation module as introduced in Sec. 3.3. Concretely, we experiment with two T-S pairs in the setting of 10% data and short schedules of 90 or 180 epochs. As shown in Table 4, PTS works as a competitive method across different setups, while scale normalization (SN) and Std Matching (SM) also

achieve performance gains upon the SVD method, which shows the importance of matching statistics while preserving original relative magnitude.

Moreover, we conduct hyper-parameter analysis to study the sensitiveness of the hyper-parameters in the PTS function. From Table 5, both T and n have a relatively wide range that could bring performance gains upon SVD (“w.o” in the Table), which illustrates the robustness of the PTS function.

T \rightarrow S	R50 \rightarrow R18	R50 \rightarrow MNV2
Epoch	90	180
SVD	75.09	76.28
SVD+SN	75.34	76.45
SVD+SM	75.87	76.89
SVD+PTS	<u>75.74</u>	77.07

Table 4. **Ablation study on the transformation module.** Setting: 10% data with 90 or 180 epochs. R50: MS. SN: Scale Normalization. SM: Std Matching. Bold: best. Underlined: second-best.

T	0.01	0.03	0.06	0.1	0.3	0.5	0.7	0.9	w.o.
Acc	75.12	75.36	75.53	75.74	75.51	75.21	75.29	75.07	75.09

n	2	3	4	5	6	w.o.
Acc	75.06	75.74	75.69	75.61	75.19	75.09

Table 5. **Hyper-parameter analysis on the PTS function.** Setting: 10% data with 90 epochs. T-S pair: MS R50 \rightarrow R18. When varying T , we fix $n = 3$. When varying n , we fix $T = 0.1$.

Ablation Studies with different teacher models. As shown in Table 6, teachers with higher accuracy on ImageNet benchmark do not necessarily lead to better KDEP performance. Intriguingly, combining the visualization results in Figure 4, we notice a strong correlation between the teacher’s feature compactness and KDEP performance, that compact feature representation suffers to serve as a good teacher, whereas diverse ones produce superior results.

Teacher	Standard SP	MS	SWSL	Swin-B	MEAL V2
ImageNet	75.77	73.85	81.12	84.81	80.68
SVD	74.29	75.09	73.75	72.053	70.33
SVD+PTS	74.79	75.74	74.27	72.252	72.41

Table 6. **Ablation study on different teacher models.** Setting: 10% data with 90 epochs. ImageNet: ImageNet val set top-1 Acc.

Ablation study on multiple layer feature-KD. Our method only distills the penultimate layer feature to transfer knowledge, whereas multiple layer feature-based KD

has also been introduced in the literature. Here, we study the effectiveness of multiple layer feature-based KD for our KDEP settings. From Table 7, we observe that multiple layer feature-based KD may be more beneficial when teacher and student are of similar architectures (*i.e.* R50 \rightarrow R18), but may cause performance degradation for dissimilar architecture pairs (*i.e.* R50 \rightarrow MNV2) due to potential semantic mismatch of intermediate features.

T \rightarrow S	MS R50 \rightarrow R18		MS R50 \rightarrow MNV2	
Epoch	90	180	90	180
SVD+PTS	75.74	77.07	75.37	76.53
SVD+PTS+ML	76.11 (+0.37)	77.21 (+0.14)	75.20 (-0.17)	76.25 (-0.28)

Table 7. **Ablation study on multiple layer feature-based KD.** Setting: 10% data with 90 or 180 epochs. ML: multiple layer feature-based KD.

Ablation study on logits KD. We also experiment with logits KD for our KDEP settings. Since the best MS R50 teacher does not contain weights to produce logits, we study with the second-best standard SP R50 teacher. As shown in Table 8, logits KD leads to inferior performance compared with our feature-based KD under various data amounts and training schedules. More importantly, logits KD may fail to utilize potential better teachers trained on more data, which in our case could produce much better KDEP results (see ‘‘Ours (SVD+PTS)[†]’’ in Table 8).

Method	Data	Epoch	Time (/h)	Acc (%)
SP. b.	10%	90	3.9	71.05
Logits	10%	90	4.0	71.95
Ours (SVD+PTS)	10%	90	4.0	74.79
Ours (SVD+PTS) [†]	10%	90	4.0	75.74
SP. b.	10%	180	7.8	71.60
Logits	10%	180	8.0	73.63
Ours (SVD+PTS)	10%	180	8.0	76.02
Ours (SVD+PTS) [†]	10%	180	8.0	77.07
SP. o.	100%	90	39	77.13
Logits	100%	90	40	77.19
Ours (SVD+PTS)	100%	90	40	77.66
Ours (SVD+PTS) [†]	100%	90	40	78.40

Table 8. **Ablation study on logits KD.** T-S pair: Standard SP R50 \rightarrow R18. [†]: MS R50 (provided weights do not contain classifiers’, so we cannot apply logits KD with this model).

4.4. Discussion

Varying KDEP data \mathcal{D}_u . In the above studies, we use 10% or 100% ImageNet-1K data as \mathcal{D}_u , which is known to have large diversity. We further experiment with different \mathcal{D}_u : (1) object-level data. We use the four downstream image classification datasets respectively; (2) scene-level data. We use COCO and ADE20K respectively. Since these different \mathcal{D}_u are of different sizes, we keep the training iterations same as the 10% data with 90 epochs setting. Results are

shown in Table 9, where we conclude that the KDEP performance could relate to different aspects of \mathcal{D}_u : data amount, data diversity, and image context (object or scene level).

Harvesting the largest data amount and diversity, ImageNet-1K yields the best KDEP performance with object-centric samples. Meanwhile, COCO holds similar data scale and fairly diverse scene-level images, producing the second-best results. On the contrary, constrained by the fewest images of only texture patterns (low diversity), DTD leads to the worst results. We suggest that enlarging data amount and diversity is beneficial for KDEP, and object-level images are more favourable for our current KDEP method. However, we believe scene-level images could be better utilized by further leveraging the characteristic of scene context, which we leave as future work.

\mathcal{D}_u	N_u	Caltech	DTD	CUB	CIFAR	Avg
Caltech	15k	72.77	67.14	67.48	80.03	71.86
DTD	3.8k	54.32	64.97	62.43	74.79	64.13
CUB	6.0k	62.74	61.72	78.82	78.31	70.40
CIFAR	50k	61.83	56.36	61.19	81.48	65.21
COCO	118k	72.34	69.55	71.68	79.87	73.36
ADE20K	20k	67.66	67.38	68.20	78.90	70.54
10% ImageNet-1K	128k	75.33	71.80	74.20	81.61	75.74

Table 9. **Varying \mathcal{D}_u .** Setting: the same training iterations as 10% ImageNet-1K data with 90 epochs. T-S pair: MS R50 \rightarrow R18.

Computation cost. In most of our experimental results, we provide the training time of KDEP and supervised pre-training, where only unnoticeable extra training time is added. Moreover, the GPU memory usage during KDEP is also similar to supervised pre-training since we do not require gradients for the teacher. Yet, our teacher model is R50, and additional computation costs may increase when using larger teachers that inquires more inference time and GPU memory. Still, the pre-training time could be largely reduced compared with supervised pre-training.

5. Conclusion

We have present KDEP, an orthogonal strategy for pre-training new models. With extensive experimental results, our simple yet efficient feature-based KD method has shown promising results for KDEP, offering several favourable characteristics: Faster Convergence, Higher Data-efficiency and Better Transferability. Without bells and whistles, KDEP achieves comparable transfer performance as supervised pre-training with only 10x less data and 5x less training time.

Limitations and Broader Impact. The KDEP performance largely relies on a suitable teacher model, where how to obtain such a teacher still needs further investigation. In our study, we found the compactness of feature distribution could be an important indicator, from which we would hope

the community could release models not only with compact features (usually tailored to ImageNet task) but also with diverse feature distributions (pre-trained with large-scale diversified data). Moreover, we hope our work could inspire more research and applications on KDEP, which we believe is very useful for academic research and practical usage.

Acknowledgement

This work has been supported in part by Hong Kong Research Grant Council - Early Career Scheme (Grant No. 27209621), HKU Startup Fund, HKU Seed Fund for Basic Research, and SmartMore donation fund.

References

- [1] Pulkit Agrawal, Ross Girshick, and Jitendra Malik. Analyzing the performance of multilayer neural networks for object recognition. In *European conference on computer vision*, pages 329–344. Springer, 2014.
- [2] 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*, pages 9163–9171, 2019.
- [3] Devansh Arpit, Stanislaw Jastrzebski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at memorization in deep networks. In *International Conference on Machine Learning*, pages 233–242. PMLR, 2017.
- [4] Cristian Bucilua, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 535–541, 2006.
- [5] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *arXiv preprint arXiv:2006.09882*, 2020.
- [6] Defang Chen, Jian-Ping Mei, Yuan Zhang, Can Wang, Zhe Wang, Yan Feng, and Chun Chen. Cross-layer distillation with semantic calibration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 7028–7036, 2021.
- [7] Hanting Chen, Yunhe Wang, Chang Xu, Chao Xu, and Dacheng Tao. Learning student networks via feature embedding. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1):25–35, 2020.
- [8] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [10] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey Hinton. Big self-supervised models are strong semi-supervised learners. *arXiv preprint arXiv:2006.10029*, 2020.
- [11] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021.
- [12] Xinyang Chen, Sinan Wang, Bo Fu, Mingsheng Long, and Jianmin Wang. Catastrophic forgetting meets negative transfer: Batch spectral shrinkage for safe transfer learning. 2019.
- [13] Ying-Cong Chen, Xiaoyong Shen, and Jiaya Jia. Makeup-go: Blind reversion of portrait edit. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4501–4509, 2017.
- [14] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3606–3613, 2014.
- [15] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016.
- [16] 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*, pages 248–255. Ieee, 2009.
- [17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [18] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015.
- [19] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.
- [20] Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long tail via influence estimation. *arXiv preprint arXiv:2008.03703*, 2020.
- [21] Gregory Griffin, Alex Holub, and Pietro Perona. Caltech-256 object category dataset. 2007.
- [22] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020.
- [23] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Con-*

- ference on Computer Vision and Pattern Recognition, pages 9729–9738, 2020.
- [24] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [27] Byeongho Heo, Jeesoo Kim, Sangdoon 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*, pages 1921–1930, 2019.
- [28] Byeongho Heo, Minsik Lee, Sangdoon Yun, and Jin Young Choi. Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3779–3787, 2019.
- [29] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [30] Zehao Huang and Naiyan Wang. Like what you like: Knowledge distill via neuron selectivity transfer. *arXiv preprint arXiv:1707.01219*, 2017.
- [31] Xiao Jin, Baoyun Peng, Yichao Wu, Yu Liu, Jiaheng Liu, Ding Liang, Junjie Yan, and Xiaolin Hu. Knowledge distillation via route constrained optimization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1345–1354, 2019.
- [32] Armand Joulin, Laurens Van Der Maaten, Allan Jabri, and Nicolas Vasilache. Learning visual features from large weakly supervised data. In *European Conference on Computer Vision*, pages 67–84. Springer, 2016.
- [33] Jangho Kim, SeongUk Park, and Nojun Kwak. Paraphrasing complex network: Network compression via factor transfer. *arXiv preprint arXiv:1802.04977*, 2018.
- [34] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pages 491–507. Springer, 2020.
- [35] Nikos Komodakis and Sergey Zagoruyko. Paying more attention to attention: improving the performance of convolutional neural networks via attention transfer. In *ICLR*, 2017.
- [36] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [37] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *IJCV*, 2020.
- [38] Seung Hyun Lee, Dae Ha Kim, and Byung Cheol Song. Self-supervised knowledge distillation using singular value decomposition. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 335–350, 2018.
- [39] Ang Li, Allan Jabri, Armand Joulin, and Laurens van der Maaten. Learning visual n-grams from web data. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4183–4192, 2017.
- [40] Xingjian Li, Haoyi Xiong, Hanchao Wang, Yuxuan Rao, Liping Liu, Zeyu Chen, and Jun Huan. Delta: Deep learning transfer using feature map with attention for convolutional networks. *arXiv preprint arXiv:1901.09229*, 2019.
- [41] Zhizhong Li, Avinash Ravichandran, Charless Fowlkes, Marzia Polito, Rahul Bhotika, and Stefano Soatto. Representation consolidation for training expert students. *arXiv preprint arXiv:2107.08039*, 2021.
- [42] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017.
- [43] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [44] Yufan Liu, Jiajiong Cao, Bing Li, Chunfeng Yuan, Weiming Hu, Yangxi Li, and Yunqiang Duan. Knowledge distillation via instance relationship graph. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7096–7104, 2019.
- [45] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. *arXiv preprint arXiv:2103.14030*, 2021.
- [46] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *Proceedings of the European conference on computer vision (ECCV)*, pages 181–196, 2018.
- [47] Aditya K Menon, Ankit Singh Rawat, Sashank Reddi, Seungyeon Kim, and Sanjiv Kumar. A statistical perspective on distillation. In *International Conference on Machine Learning*, pages 7632–7642. PMLR, 2021.
- [48] Rafael Müller, Simon Kornblith, and Geoffrey Hinton. When does label smoothing help? *arXiv preprint arXiv:1906.02629*, 2019.
- [49] 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*, pages 3967–3976, 2019.
- [50] Nikolaos Passalis and Anastasios Tefas. Learning deep representations with probabilistic knowledge transfer. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 268–284, 2018.

- [51] Nikolaos Passalis, Maria Tzelepi, and Anastasios Tefas. Heterogeneous knowledge distillation using information flow modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2339–2348, 2020.
- [52] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32:8026–8037, 2019.
- [53] Baoyun Peng, Xiao Jin, Jiaheng Liu, Dongsheng Li, Yichao Wu, Yu Liu, Shunfeng Zhou, and Zhaoning Zhang. Correlation congruence for knowledge distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5007–5016, 2019.
- [54] Hieu Pham, Zihang Dai, Qizhe Xie, and Quoc V Le. Meta pseudo labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11557–11568, 2021.
- [55] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [56] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28:91–99, 2015.
- [57] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lih Zelnik-Manor. Imagenet-21k pretraining for the masses. *arXiv preprint arXiv:2104.10972*, 2021.
- [58] 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.
- [59] Hadi Salman, Andrew Ilyas, Logan Engstrom, Ashish Kapoor, and Aleksander Madry. Do adversarially robust imagenet models transfer better? *arXiv preprint arXiv:2007.08489*, 2020.
- [60] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018.
- [61] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *Proceedings of the IEEE international conference on computer vision*, pages 843–852, 2017.
- [62] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. *arXiv preprint arXiv:1910.10699*, 2019.
- [63] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1365–1374, 2019.
- [64] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [65] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- [66] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10687–10698, 2020.
- [67] LI Xuhong, Yves Grandvalet, and Franck Davoine. Explicit inductive bias for transfer learning with convolutional networks. In *International Conference on Machine Learning*, pages 2825–2834. PMLR, 2018.
- [68] Ze Yang, Tiange Luo, Dong Wang, Zhiqiang Hu, Jun Gao, and Liwei Wang. Learning to navigate for fine-grained classification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 420–435, 2018.
- [69] Yazhou Yao, Jian Zhang, Xian-Sheng Hua, Fumin Shen, and Zhenmin Tang. Extracting visual knowledge from the internet: making sense of image data. In *International Conference on Multimedia Modeling*, pages 862–873. Springer, 2016.
- [70] Mang Ye, Xu Zhang, Pong C Yuen, and Shih-Fu Chang. Un-supervised embedding learning via invariant and spreading instance feature. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6210–6219, 2019.
- [71] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4133–4141, 2017.
- [72] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? *arXiv preprint arXiv:1411.1792*, 2014.
- [73] Kaichao You, Zhi Kou, Mingsheng Long, and Jianmin Wang. Co-tuning for transfer learning. *Advances in Neural Information Processing Systems*, 33, 2020.
- [74] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. *arXiv preprint arXiv:2103.03230*, 2021.
- [75] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017.
- [76] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127(3):302–321, 2019.
- [77] Helong Zhou, Liangchen Song, Jiajie Chen, Ye Zhou, Guoli Wang, Junsong Yuan, and Qian Zhang. Rethinking soft labels for knowledge distillation: A bias-variance tradeoff perspective. *arXiv preprint arXiv:2102.00650*, 2021.

Outline

In this supplementary file, we first provide more results in Sec. A: results of MobileNetV2 as student in Sec. A.1, detailed proof of Theorem 1 in Sec. A.2, comparison of previous feature-based KD methods in Sec. A.3, and an example of local transformations breaking the original relative magnitude in Sec. A.4. Further, we offer the elaborated implementation details for the KDEP setups and downstream task setups in Sec. B.

A. More Results

A.1. Main Results: MobileNetV2 as Student.

Due to the length limit of the main paper, we show the results of MobileNetV2 as student in Table S.1. Similar results have been achieved with MobileNetV2 (MNV2) as student compared to R18 as student, which shows the generalization of the proposed KDEP method.

A.2. Detailed Proof of Theorem 1

Given two independent random variables with normal distribution $T \sim N(0, \sigma^2)$ and $S \sim N(0, \sigma_s^2)$, then $F(\sigma) = \mathbb{E}[(T - S)^2]$ is monotonically increasing ($\sigma > 0$).

Proof 2 (Detailed version)

$$\begin{aligned} F(\sigma) &= \mathbb{E}[(T - S)^2] \\ &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} (t - s)^2 \cdot P(t, s) dt ds \\ &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} (t - s)^2 \cdot P(t) \cdot P(s) dt ds \\ &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} (t - s)^2 \cdot \frac{1}{2\pi\sigma\sigma_s} e^{-\frac{s^2}{2\sigma_s^2}} e^{-\frac{t^2}{2\sigma^2}} dt ds \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} e^{-\frac{t^2}{2\sigma^2}} \left(\int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}\sigma_s} (t - s)^2 \cdot e^{-\frac{s^2}{2\sigma_s^2}} ds \right) dt \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} e^{-\frac{t^2}{2\sigma^2}} (t^2 + \sigma_s^2) dt = \sigma^2 + \sigma_s^2 \\ \frac{dF(\sigma)}{d\sigma} &= 2\sigma > 0 \Rightarrow \text{monotonically increasing} \end{aligned}$$

A.3. Compare other feature-based KD methods.

Here, we show the KDEP results of some traditional feature-based KD methods that are developed for distilling knowledge to improve student’s performance for a specific task instead of its transferability. Also, these methods all require task label loss which violates our setting of an unlabeled dataset. Hence, we don’t include them in our paper to compare for fairness. As shown in Table S.2, previous feature-based KD methods largely rely on logit KD loss and task label loss, and perform inferiorly with only feature-level clues.

Method	Data	Epoch	Time (h)	Classification (Acc %)				
				Caltech	DTD	CUB	CIFAR	Avg
rand. init.	-	-	-	51.77	57.34	60.44	76.66	61.55
SP. b.	10%	90	4.2	66.58	65.72	71.09	78.60	70.50
KDEP	10%	90	4.3	74.34	71.84	74.24	81.06	75.37
SP. b.	100%	9	4.2	68.09	67.514	73.33	78.95	71.97
KDEP	100%	9	4.3	74.48	72.51	74.76	81.47	75.81
SP. b.	10%	180	8.4	69.23	67.33	73.62	79.50	72.42
KDEP	10%	180	8.6	75.56	73.29	75.28	81.98	76.53
SP. b.	100%	18	8.4	71.83	69.60	74.64	80.13	74.05
KDEP	100%	18	8.6	76.06	73.14	76.00	82.15	76.83
SP. b.	10%	900	42	69.93	67.59	72.74	79.83	72.52
KDEP	10%	900	43	77.66	73.05	76.06	82.53	77.32
SP. o.	100%	90	42	76.43	72.26	76.34	81.91	76.74
KDEP	100%	90	43	78.57	73.94	76.40	82.73	77.91

Method	Data	Epoch	Time (h)	Segmentation (mIoU %)			
				Cityscapes	VOC12	ADE20K	Avg
rand. init.	-	-	-	40.33	39.23	23.07	34.21
SP. b.	10%	90	4.2	60.92	62.60	29.17	50.89
KDEP	10%	90	4.3	63.50	67.28	31.46	54.08
SP. b.	100%	9	4.2	61.42	64.31	29.61	51.78
KDEP	100%	9	4.3	63.53	68.17	32.00	54.57
SP. b.	10%	180	8.4	61.68	64.65	29.95	52.09
KDEP	10%	180	8.6	64.32	68.73	31.92	54.99
SP. b.	100%	18	8.4	62.23	65.82	29.55	52.53
KDEP	100%	18	8.6	64.23	69.32	32.41	55.32
SP. b.	10%	900	42	61.87	64.95	31.07	52.63
KDEP	10%	900	43	63.89	70.45	32.23	55.52
SP. o.	100%	90	42	64.16	69.48	31.69	55.11
KDEP	100%	90	43	64.72	71.07	32.39	56.06

Method	Data	Epoch	Time (h)	Detection					
				VOC0712			COCO		
				AP	AP ₅₀	AP ₇₅	AP	AP ₅₀	AP ₇₅
rand. init.	-	-	-	31.4	56.7	30.2	25.8	43.6	26.9
SP. b.	10%	90	4.2	42.5	71.3	44.0	27.7	46.4	29.0
KDEP	10%	90	4.3	46.7	75.7	49.0	29.7	48.8	31.2
SP. b.	100%	9	4.2	43.0	71.9	44.0	27.9	46.4	29.0
KDEP	100%	9	4.3	47.1	76.3	49.6	30.2	49.8	31.9
SP. b.	10%	180	8.4	43.0	71.4	45.0	27.6	46.2	29.0
KDEP	10%	180	8.6	47.0	75.9	49.9	30.0	49.4	31.5
SP. b.	100%	18	8.4	43.7	72.6	45.4	27.9	46.6	29.3
KDEP	100%	18	8.6	47.2	76.3	50.0	30.4	50.1	32.0
SP. b.	10%	900	42	44.0	73.1	45.6	28.7	47.9	30.1
KDEP	10%	900	43	47.0	76.0	49.8	29.7	49.2	31.4
SP. o.	100%	90	42	45.5	75.0	47.6	29.6	49.1	31.3
KDEP	100%	90	43	46.8	76.5	49.4	30.0	49.6	31.4

Table S.1. KDEP vs. SP, R50 → MNV2, fine-tuned on various tasks. KDEP refers to our SVD+PTS method.

Method	SP	FitNet [58]	AT [35]	NST [30]	AB [28]	Heo [27]	Ours
Acc	71.05	72.43	67.84	67.05	71.66	72.98	75.74

Table S.2. Compare feature-based KD with 10% data and 90 epochs. Acc: averaged top-1 accuracy over 4 classification tasks.

A.4. Example: Breaking Relative Magnitude.

In Sec. 3.3 of our paper, we argue that Scale Normalization (SN) and Std Matching (SM) are local transformations that transform channel-wisely, which may break the original relative magnitude between channels. Here, we provide a toy example as an illustration.

For instance, we have a three channel penultimate layer

with target Std=[4, 3, 2] and after SVD Std=[50, 5, 1]. For a feature after SVD that is [10, 2, 2], with SN we have [0.2, 0.4, 2], and with SM we have [0.8, 1.2, 4], both losing the original relative magnitude.

B. Implementation Details

We implement our method using the PyTorch [52] framework and use SGD with momentum of 0.9 for all our experiments. All experiments are conducted using four 32G V100 GPUs.

B.1. KDEP Setups

For the KDEP procedure, we use an initial learning rate (lr) of 0.3 for R50→R18 and 0.1 for R50→MNv2. Batch size is set to 512. For data augmentation, we use RandomResizedCrop(224) and RandomHorizontalFlip. In order to reduce the burden of hyper-parameter tuning (*e.g.* weight decay), we multiply our feature-based loss (refer to Eq. 1 in our paper) by a loss weight w , which matches the feature-based loss to the loss scale of supervised pretraining. For reproduction, we provide the loss weight of different teacher-student pairs in Table S.3.

Teacher→Student	w
Standard SP R50→R18	20
MS R50→R18	3
SWSL R50→R18	1
MEAL V2 R50→R18	1
Swin-B→R18	3
MS R50→MNv2	3

Table S.3. Value of loss weight w for different T-S pairs.

For the 10% ImageNet data setting, we sample 10% images from each class of the original 1000 class in ImageNet-1K. We set the training epochs to 90 or 180, and drop the lr by a factor of 10 at 1/3 and 2/3 of total epochs. When using all of ImageNet data, we train for 9 or 18 epochs to verify fast convergence, and for 90 epochs to further boost performance, where 90 epochs with all ImageNet data is the standard supervised pretraining schedule [59]. We use weight decay $\in \{1e-4, 4e-4, 5e-4\}$ according to the length of the training schedule, shown in Table S.4.

Data	Epoch	Weight decay
10%	90	5e-4
10%	180	4e-4
100%	9	5e-4
100%	18	4e-4
100%	90	1e-4

Table S.4. Weight decay for different KDEP training scedules.

B.2. Downstream Task Setups

For all downstream tasks, we use the same schedule and evaluation protocols for all models for a fair comparison.

For image classification, we initialize the backbone with the distilled weights and add a linear classifier with random initialization. We train the network for 150 epochs with batch size of 64, weight decay of $5e-4$, an initial learning rate $\in \{0.01, 0.001\}$ which drops by a factor of 10 at 1/3 and 2/3 of total epochs. Again, we use RandomResized-Crop(224) and RandomHorizontalFlip for data augmentation.

For semantic segmentation, we also initialize the backbone with the distilled weights, and add a PSP module [75] and a segmentation head after the backbone. We use batch size of 16, an initial learning rate of 0.01, weight decay of $1e-4$, crop size of 512, and deploy an polynomial learning rate annealing procedure [8]. For data augmentation, we use random scaling, random horizontal flipping, random rotation, and random Gaussian blur. The number of epochs is 50, 100, 200 for VOC12, ADE20K, and Cityscapes, respectively, following previous standard [75].

For object detection, we experiment with the Faster R-CNN [56] detector and backbones are also initialized by the distilled weights. Unless noted, all the setups follow the evaluation protocols in MOCO [23]. For ResNet18, we use a backbone of R18-C4 (similar to R50-C4 [24]) for both VOC and COCO experiments. For MobileNetV2, we equip it with a FPN [42] backbone. 1x schedule is applied for COCO.

B.3. Parametric/Non-parametric Methods.

In our experiments, we experiment with three 1×1 conv variants for parametric methods: Post-ReLU, Pre-ReLU and the one in [27]. Specifically, we add a 1×1 convolutional layer and a Batch Normalization layer either after (Post-ReLU) or before (Pre-ReLU) the ReLU activation function. We also experiment with the Pre-ReLU 1×1 conv method equipped with Margin ReLU of teacher’s feature and Partial L_2 loss function as in [27].

For non-parametric methods, we explore channel selection (CS.var, CS.rand), interpolation, and SVD. For channel selection methods, we experiment with two methods: selecting the top- D_s channels with largest variances (CS.var) or random selecting D_s channels (CS.rand). For interpolation method, we use the default nearest-neighbor interpolation in PyTorch [52]. For SVD, we calculate the singular vectors offline and use the top- D_s principal components to transform the teacher’s features during the online KDEP process.