

DistPro: Searching A Fast Knowledge Distillation Process via Meta Optimization

Xueqing Deng^{*1,3}, Dawei Sun^{*2}, Shawn Newsam³, and Peng Wang¹

¹ ByteDance Inc. {xueqingdeng,peng.wang}@bytedance.com

² ECE, UIUC, {daweis2}@illinois.edu

³ EECS, UC Merced, {snewsam}@ucmerced.edu

Abstract. Recent Knowledge distillation (KD) studies show that different manually designed schemes impact the learned results significantly. Yet, in KD, automatically searching an optimal distillation scheme has not yet been well explored. In this paper, we propose **DistPro**, a novel framework which searches for an optimal KD process via differentiable meta-learning. Specifically, given a pair of student and teacher networks, **DistPro** first sets up a rich set of KD connections from the transmitting layers of the teacher to the receiving layers of the student, and in the meanwhile, various transforms are also proposed for comparing feature maps along their pathways for distillation. Then, each combination of connection and transform (pathway) is associated with a stochastic weighting process which indicates its importance at every step during the distillation. At the searching stage, the process can be effectively learned through our proposed bi-level meta-optimization strategy. At the distillation stage, **DistPro** adopts the learned processes for knowledge distillation, which significantly improves the student accuracy especially when faster training is required. Lastly, we find the learned processes can be generalized between similar tasks and networks. In our experiments, **DistPro** produces state-of-the-art (SoTA) accuracy under varying number of learning epochs on popular datasets, *i.e.* CIFAR100 and ImageNet, which demonstrates the effectiveness of our framework. Codes are available at <https://github.com/xdeng7/DistPro>.

1 Introduction

Knowledge distillation (KD) is proposed to effectively transfer knowledge from a well performing larger/teacher deep neural network (DNN) to a given smaller/student network, where the learned student network often learns faster or performs better than that with a vanilla training strategy solely using ground truths.

Since its first appearance in DNN learning [14], KD has achieved remarkable success in training efficient models for image classification [50], image segmentation [26], object detection [3], *etc*, contributing to its wide application in various model deployment over mobile-phones or other low-power computing devices [28]. Nowadays, KD has become a popular technique in industry to develop distilled DNNs to deal with billions of data per day.

* These authors contributed equally to this work

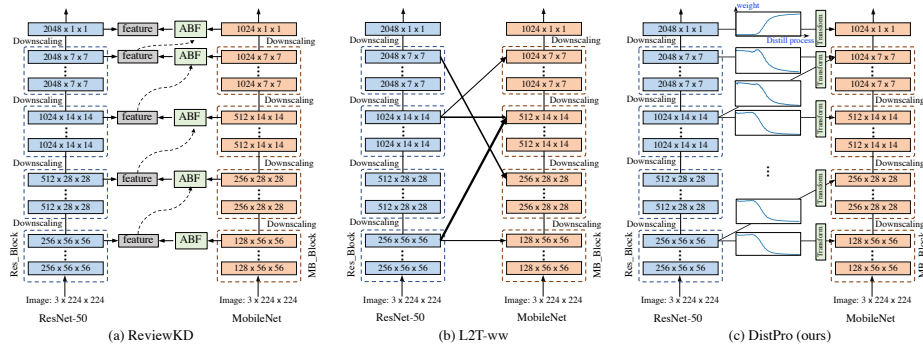


Fig. 1: Comparisons of distillation methods for learning process between teacher and student models. (a) Knowledge Review [4] proposes fixed sampled pathways by enumerating different configurations. (b) L2T-ww [18], adopts a meta-learning framework to learn a floating weight for each selected pathway. (c) Our framework learns distillation process for each pathway.

To improve the distillation efficiency and accuracy, numerous handcrafted KD design schemes have been proposed, *e.g.*, designing different distillation losses at outputs [26, 44], manually assigning intermediate features maps for additional KD guidance [50, 49]. However, recent studies [11, 27, 48] indicate that the effectiveness of those proposed KD techniques varies on the networks and tasks. Some recent works propose to search some configurations to conclude a better KD scheme. For instance, ReviewKD [4] (Fig. 1(a)) proposes to evaluate a subset from a total of 16 pathways by enabling and disabling partial of them for image classification task. It comes to the conclusion that partial pathways are always redundant. However, our experimental results show that it does not hold to a semantic segmentation task, and after conducting a research of the pathways, we may obtain better results. L2T-ww [18] (Fig. 1(b)) takes a further step, it can not only set up multiple distillation pathways between feature maps, but also learn a floating weight for each pathway, showing better performance than fixed weight. This inspires us to explore better KD schemes, and as motivated by the learning rate scheduler, shown in Fig. 1(c), we propose the concept of distillation process for each pathway i , *i.e.* $\mathcal{A}^i = \{\alpha_t^i\}_{t=1}^T$, where T is the number of distillation steps. Thus, it can result in dynamic changes of the importance for each pathway along the whole distillation procedure. Our work shows that such a process is beneficial to acceleration.

However, searching a process is more difficult than finding a floating weight, which includes T times more parameters. Obviously, it is not practical to solve the goal via a brute force way, for example, randomly drawing a sample process, and later fully training and validating the network to evaluate its performance. Thanks to bi-level meta-learning [10, 25], we find our problem can be formulated and tackled in a similar manner. Such a framework not only skips the difficulty of random exploration through valid meta gradient, but also naturally provides soft weighting that can be adopted to generate the process. Additionally, to effectively

apply the framework and avoid possible noisy gradient from the meta-training, we propose a proper normalization for each $\alpha_t = [\alpha_t^0, \dots, \alpha_t^N]$, where N is the number of pathways. We call our framework **DistPro** (Distillation Process). In our experiments, we show that **DistPro** produces better results on various tasks, such as classification with CIFAR100 [22] and ImageNet1K [7], segmentation with CityScapes [6] and depth estimation with NYUv2 [32].

Finally, we find our learned process remains similar with minor variations across different network architectures and tasks as long as it uses the same proposed pathways and transforms. This indicates that the process can be generalized to new tasks. In practice, we transfer the process learned by CIFAR100 to ImageNet1K, and show that it improves over the baselines and accelerates the distillation (2x faster than ReviewKD [4] as shown in Tab. 3).

In summary, our contributions are three-folds 1) We propose a meta-learning framework for KD, *i.e.* DistPro, to efficiently learn an optimal process to perform KD. 2) We verify **DistPro** over various configurations, architecture and task settings, yielding significant improvement over other SoTA methods. 3) Through the experiments, we find the process that can generalize across tasks and networks, can potentially benefit KD in new tasks without additional searching.

2 Related Works

Knowledge distillation. Starting under the name knowledge transfer [2, 42], knowledge distillation (KD) is later popularized owing to Hinton et.al [14] for training efficient neural networks. Thereafter, it has been a popular field in the past few years, in terms of designing KD losses [43, 44], combination with multiple tasks [33, 8] or dealing with specific issues, eg few-shot learning [20], long-tail recognition [45]. Here, we majorly highlight the works that are closely related to ours, in order to locate our contributions.

According to a recent survey [11], current KD literature includes multiple knowledge types, eg response-based [31], feature-based [35, 13] and relation-based [41] knowledge. For distillation algorithms, various distillation criteria are proposed such as adversarial-based [43], attention-based [17], graph-based [24] and lifelong distillation [5] etc. Finally, based on a certain task, KD can also extend with different task-aware metrics, eg speech [33], NLP [8] etc. In our principle, we hope all the surveyed KD schemes, ie knowledge types, methods under certain task settings, can be pooled with a universal way to a search space, in order to find the best distillation. While in this paper, we take the first step towards this goal by exploring a sub-field in this whole space, which is already a challenging problem to solve. Specifically, we adopt the setting of offline distillation with feature-based and response-based knowledge, where both network responses and intermediate feature maps are adopted for KD. For KD method, we use attention-based methods to compare feature responses, and apply the KD model to vision tasks including classification, segmentation and depth estimation.

Inside this field, knowledge review [4] and L2T-ww [18] are the most related to our work. The former investigates the importance of a few pathways and propose a knowledge review mechanism with a novel connection pattern, ie, residual learning framework. It provides SoTA results in several commonly comparison benchmarks. The latter learns a fixed weights for a few intermediate pathways for few-shot knowledge transfer. As shown in Fig. 1, **DistPro** finds a distillation process. Therefore, we extend the search space. In addition, for dense prediction tasks, one related work is IFVD [44], which proposes an intra-class feature variation comparison (IFVD). **DistPro** is free to extend to dense prediction tasks, and it also obtain extra benefits after combined with IFVD.

Meta-learning for KD/hyperparameters. To automate the learning of a KD scheme, we investigated a wide range of efficient meta-learning methods in other fields that we might adopt. For examples, L2L [1] proposes to learn a hyperparameter scheduling scheme through a RNN-based meta-network. Franceschi et.al [10] propose an gradient-based approach without external meta-networks. Later, these meta-learning ideas have also been utilized in tasks of few-shot learning (eg learning to reweight [34]), learning cross-task knowledge transfer(eg learning to transfer [18]), and neural architecture search (NAS)(eg DARTS [25]). Although these methods share similar framework, it is critical to have essential embedded domain knowledge and task-aware adjustment to make it work. In our case, inspired by these methods, we majorly utilize the gradient-based strategy due to its efficiency for KD scheme learning, and also first to propose using the learnt process additional to the learnt importance factor.

Finally, knowledge distillation for NAS has also drawn significant attention recently. For example, Liu et.al [27] try to find student models that are best for distilling a given teacher, while Yao et.al [48] propose to search architectures of both the student and the teacher models based on a certain distillation loss. Though different from our scenario, *i.e.* fixed student-teacher architectures, it raises another important question of how to find the Pareto optimal inside the union space of architectures and KD schemes under certain resource constraint, which we hope could inspire future researches.

3 Approach

In this section, we elaborate **DistPro** by first setting up KD pathways with intermediate features in order to construct our search space, establishing the notations and definition of our KD scheme. Then, we derive the gradient for generating our process for the scheme. At last, the overall algorithm is presented.

3.1 KD with intermediate features

Numerous prior works [19, 4] have demonstrated that intermediate feature maps from neural network can benefit distillation. Motivated by this, we design our approach using feature maps. Let the student neural network be \mathcal{S} and the

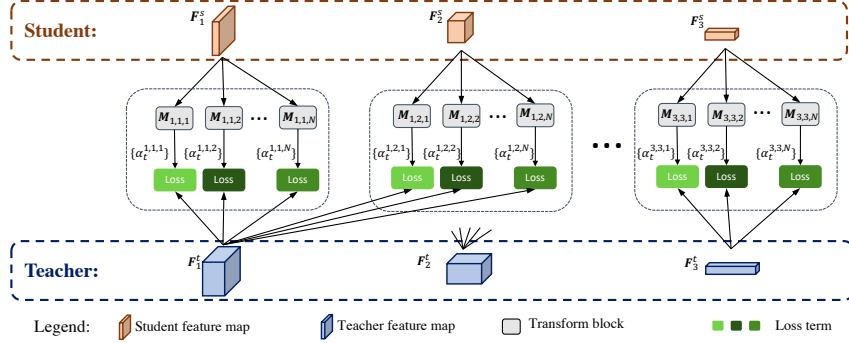


Fig. 2: Illustration of the search space. Two groups of feature maps are selected from the student and the teacher. For each pair of feature maps \mathcal{F}_j^s and \mathcal{F}_i^t , we insert N candidate transform blocks $M_{i,j,1}, M_{i,j,2}, \dots, M_{i,j,N}$ and get N loss terms. For each loss term, a distill process $\mathcal{A}^{i,j,k} = \{\alpha_t^{i,j,k}\}_{t=1}^T$ is assigned.

teacher neural network be \mathcal{T} . Given an input example \mathbf{X} , the output of the student network is as follows

$$\mathcal{S}(\mathbf{X}) := \mathcal{S}_{L_s} \circ \dots \circ \mathcal{S}_2 \circ \mathcal{S}_1(\mathbf{X}), \quad (1)$$

where \mathcal{S}_i is the i -th layer of the neural network and L_s is the number of layers. The k -th intermediate feature map \mathcal{F}_k^s of the student is defined as follows,

$$\mathcal{F}_k^s(\mathbf{X}) := \mathcal{S}_k \circ \dots \circ \mathcal{S}_2 \circ \mathcal{S}_1(\mathbf{X}), \quad 1 \leq k \leq L_s.$$

Similarly, the feature map of the teacher is denoted by \mathcal{F}_k^t , $1 \leq k \leq L_t$.

Knowledge can be distilled from a pathway between i -th feature map of the teacher and j -th feature map of the student by penalizing the difference between these two feature maps, ie, \mathcal{F}_i^t and \mathcal{F}_j^s . Since the feature maps may come from any stage of the network, they may not be in the same shape and thus not directly comparable. Therefore, additional computation are required to align these two feature maps into the same shape. To this end, a transform block is added after \mathcal{F}_j^s , which could be in any forms of differentiable computation. In our experiments, the transform block consists of multiple convolution layers and an interpolation layer to align the spatial and channel resolution of the feature map. Denoting the transform block by \mathcal{M} , then the loss term measuring the difference between these two feature maps is as follows,

$$\ell(\mathcal{F}_j^s, \mathcal{F}_i^t) := \delta(\mathcal{M}(\mathcal{F}_j^s), \mathcal{F}_i^t),$$

where δ is the optional distance function, which could be L1 distance or L2 distance as used in [19], etc.

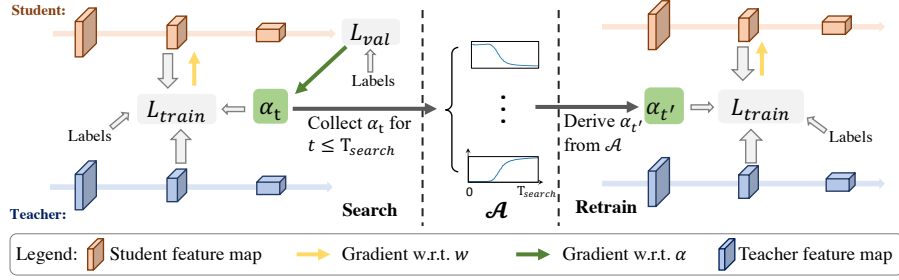


Fig. 3: Two phases of the proposed algorithm. In the search phase, student θ and α_t are computed. In the retrain phase, the learned process of $\mathcal{A} = \{\alpha_t, 0 \leq t \leq T_{search}\}$ is interpolated to be used for KD, and only θ is updated.

3.2 The Distillation Process

Now we will be able to build pathways to connect any feature layer of the teacher to any layer of the student with appropriate transforms. However, as discussed in Sec. 1, not all of these pathways are beneficial. This motivates us to design an approach to find out the importance of each possible pathway by assigning an importance factor for it. Different from the existing work [18, 4], the importance factor here is a process. Formally, a stochastic weight process $\mathcal{A}^i = \{\alpha_t^i\}_{t=1}^T$ is associated with the pathway i , where T is the total learning steps. Here, α_t^i describes the importance factor at different learning step t .

Let $D_{train} := \{(\mathbf{X}_i, y_i)\}_{i=1}^{|D_{train}|}$ and $D_{val} := \{(\mathbf{X}_i, y_i)\}_{i=1}^{|D_{val}|}$ be the training set and the validation set respectively, where y_i is the label of sample \mathbf{X}_i . We assume that for each pair of feature maps, \mathcal{F}_i^t and \mathcal{F}_j^s , we have N candidate pre-defined transforms, $\mathbf{M}_{i,j,1}, \mathbf{M}_{i,j,2}, \dots, \mathbf{M}_{i,j,N}$. The connections together with the transforms construct our search space as shown in Fig. 2.

We now define the search objective which consists of the optimizations of the student network and \mathcal{A} . The student network is trained on the training set with a loss encoding the supervision from both the ground truth label and the neural network. Specifically, denoting the parameters of the student and the transforms by θ , the loss on the training set at learning step t is defined as follows,

$$\begin{aligned}
 L_{train}(\theta_t, \alpha_t) &= \frac{1}{|D_{train}|} \sum_{(\mathbf{X}, y) \in D_{train}} \left(\delta_{label}(\mathcal{S}(\mathbf{X}), y) \right. \\
 &\quad \left. + \sum_{i=1}^{L_t} \sum_{j=1}^{L_s} \sum_{k=1}^N \alpha_t^{i,j,k} \delta(\mathbf{M}_{i,j,k}(\mathcal{F}_j^s(\mathbf{X})), \mathcal{F}_i^t(\mathbf{X})) \right), \\
 \alpha_t^{i,j,k} &= \frac{\exp(\tilde{\alpha}_t^{i,j,k})}{\sum_{i,j,k} \exp(\tilde{\alpha}_t^{i,j,k}) + \exp(g)}
 \end{aligned} \tag{2}$$

where $\alpha_t^{i,j,k} \in \mathbb{R}_{\geq 0}^{L_{te} \times L_s \times N}$ are the importance factors at training step t and δ_{label} is a distance function measuring difference between predictions and labels. Here, $\alpha_t = [\alpha_t^{0,0,0}, \dots]$ controls the importance of all knowledge distill pathway at training step t . For numerical stability and avoiding noisy gradients, we apply a biased softmax normalization with parameter $g = 1$ to compute α_t which is validated from various normalization strategies in our experiments. This is a commonly adopted trick in meta-learning for various tasks such as NAS [46, 25] or few-shot learning [34, 18]. More details about the discussion of normalization methods can be found in the experimental section (Sec. 4.4).

Next, our goal is to find an optimal sampled process \mathcal{A}^* yielding best KD results, where the validation set is often used to evaluate the performance of the student trained on unseen inputs. To this end, following [25], the validation loss is adopted to evaluate the quality of \mathcal{A} , which is defined as follows,

$$L_{val}(\theta) := \frac{1}{|D_{val}|} \sum_{(\mathbf{x}, y) \in D_{val}} \delta_{label}(\mathcal{S}(\mathbf{X}), y).$$

Finally, we formulate the bi-level optimization problem over \mathcal{A} and network parameter θ as:

$$\begin{aligned} \min_{\mathcal{A}} \quad & L_{val}(\theta^*(\mathcal{A})) \\ \text{s.t.} \quad & \theta^*(\mathcal{A}) = \arg \text{f}_{\theta} L_{train}(\Theta, \mathcal{A}). \end{aligned} \quad (3)$$

where $\theta^*(\mathcal{A})$ is the parameters of the student neural network trained with the loss process defined with \mathcal{A} , *i.e.* $L_{train}(\Theta, \mathcal{A}) = \{L_{train}(\theta_t, \alpha_t)\}_{t=1}^T$. The ultimate goal of the optimization is to find an \mathcal{A} such that the loss on the validation set is minimized. Note here, similar problem has been proposed in NAS [46, 25] asking for a fixed architecture, but our problem is harder, and can be optimized if we enforce all the value in \mathcal{A} to be the same. However, similarly, we are able to apply gradient-based method following the chain-rule to solve this problem.

3.3 Learning the Process \mathcal{A}

From the formulation in Eq. 3, directly solving the issue is intractable. Therefore, we propose two assumptions to simplify the problem. First, smooth assumption which means the next step of the process should be closed to previous one. This is commonly adopted in DNN training with stochastic gradient decent (SGD) with a learning rate scheduler [21]. Second, similar learning procedure assumption, which means among different distillation training given a teacher/student pair, at same step t , the student will have similar parameter θ_t . This assumption allows us to search for \mathcal{A} with a single training procedure, which also holds from our experiments when batch size is relatively large (*e.g.* 512 for ImageNet).

Based on these assumptions, we let $\alpha_{t+1} = \alpha_t + \gamma \Delta(\alpha_t | \theta_t)$, where γ controls its changing ratio to be small. Then, we are able to adopt a greedy strategy to break the original problem of Eq. 3 down to a sequence of single steps of optimization, which can be defined as,

$$\begin{aligned} \alpha_{t+1} &= \alpha_t - \gamma \nabla_{\alpha} L_{val}(\theta_{t+1}(\alpha_t)) \\ \text{s.t.} \quad \theta_{t+1}(\alpha_t) &= \theta_t - \xi \nabla_{\theta} L_{train}(\theta_t, \alpha_t) \end{aligned}$$

Algorithm 1 DistPro

Input: Full train data set D ; Pre-trained teacher; Initialization of α_0, θ_0 ; Number of iterations T_{search} and T .

Output: Trained student.

- 1: Split D into D_{train} and D_{val} .
 - 2: Let $t = 0$.
 - 3: **while** $t < T_{search}$ **do**
 - 4: Compute α_t with descending the gradient approximation in Eq. 5 and do normalization.
 - 5: Update θ_t by descending $\nabla_{\theta} L_{train}(\theta_t; \alpha_t)$ in Eq. 2.
 - 6: Push the current α_t to \mathcal{A} .
 - 7: **end while**
 - 8: Interpolate \mathcal{A} with length of T .
 - 9: Set $D_{train} = D$ or use a new D_{train} in another task.
 - 10: Reset $t = 0$.
 - 11: **while** $t < T$ **do**
 - 12: Load the α_t corresponding to the current t .
 - 13: Update θ by descending $\nabla_{\theta} L_{train}(\theta, \alpha_t)$.
 - 14: **end while**
-

where ξ is the learning rate of inner optimization for training student network. In practice, the inner optimization can be solved with more sophisticated gradient-based method, eg, gradient descent with momentum. In those cases, Eq. 4 has to be modified accordingly, but all the following analysis still applies.

Next, we apply the chain rule to Eq. 4 and get

$$\nabla_{\alpha} L_{val}(\theta_t - \xi \nabla_{\theta} L_{train}(\theta_t, \alpha_t)) = -\xi \nabla_{\alpha, \theta}^2 L_{train}(\theta_t, \alpha_t) \nabla_{\theta} L_{val}(\theta_{t+1}), \quad (4)$$

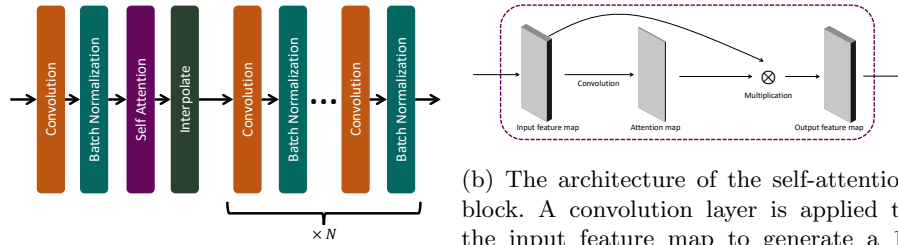
However, the above expression contains second-order derivatives, which is still computational expensive. Next, we approximate this second-order derivative with finite difference as introduced in [25]. Let ϵ be a small positive scalar and define the notation $\theta^{\pm} = \theta \pm \epsilon \nabla_{\theta} L_{val}(\theta_{t+1})$. Then, we have

$$\nabla_{\alpha, \theta}^2 L_{train}(\theta_t, \alpha_t) \nabla_{\theta} L_{val}(\theta_{t+1}) = \frac{\nabla_{\alpha} L_{train}(\theta^+, \alpha_t) - \nabla_{\alpha} L_{train}(\theta^-, \alpha_t)}{2\epsilon}. \quad (5)$$

Finally, we set an initial value $\alpha_0 = 1$ and launch the greedy learning procedure. Once learned, we push all computed α_t with T steps to a sequence, which is served as a sample of learned stochastic distillation process \mathcal{A} , and it can be used for retraining the student network with KD in the same or similar tasks.

3.4 Acceleration and adopting \mathcal{A} for KD.

Now let us take a closer look at the approximation. To evaluate the expression in Eq. 5, the following items have to be computed. First, computing θ_{t+1} requires a forward and backward pass of the student, and a forward pass of the teacher. Then, computing θ^{\pm} requires another forward and backward pass of the student.



(a) The architecture of the transform block. The self-attention block is illustrated in Fig. 4b.

(b) The architecture of the self-attention block. A convolution layer is applied to the input feature map to generate a 1-channel attention map. Then the input feature map is multiplied with the attention map.

Fig. 4: Transform blocks

Finally, computing $\nabla_{\alpha} L_{train}(\theta^{\pm}, \alpha_t)$ requires two forward passes of the student. In conclusion, evaluating the approximated gradient in Eq. 5 entails one forward pass of the teacher, and four forward passes and two backward passes of the student in total. This is a time consuming process, especially when the required KD learning epoch is large, *e.g.* ≥ 100 .

In meta-learning NAS literature, to avoid 2_{nd} order approximation, researchers commonly adopt 1_{st} order training [46, 25] solely based on training loss. However, this is not practical in our case since the hyperparameters are defined over training loss. 1_{st} order gradient over α_t will simply drive all value to 0. Therefore, in our case, we choose to reduce the learning epochs of \mathcal{A} to T_{search} , which is much smaller than T , and then expand it to a sequence with length T using linear interpolation. The choices of T_{search} can be dynamically adjusted based on the dataset size, which will be elaborated in Sec. 4.

Additionally, since we have more pathways than other KD strategy [4], therefore, the loss computation cost can not be ignored. To reduce the cost, at step t , we use a clip function to all $\tilde{\alpha}_t$ in Eq. 2 with a threshold $\tau = 0.5$, and drop corresponding computation when $\tilde{\alpha}_t \leq \tau$. This can save us 60% of loss computational cost in average, resulting in comparable KD time with our baselines.

In summary, the overall DistPro process is presented in Alg. 1, where two phases of the algorithm are explained in order. The first phase is for searching the scheme \mathcal{A} . To this end, α_t and θ_t are computed alternately. The α_t obtained in each step is stored for future usage. The Second phase is for retraining the neural network with all the available training data and the searched \mathcal{A} .

4 Experiments

In this section, we evaluate the proposed approach on several benchmark tasks including image classification, semantic segmentation, and depth estimation. For image classification, we consider the popularly used dataset CIFAR100 and ImageNet1K. For semantic segmentation and depth estimation, we consider CityScapes [6] and NYUv2 [32] respectively. To make fair comparison, we use

Teacher	WRN-40-2	WRN-40-2	ResNet32x4	ResNet32x4	ResNet56	ResNet110
Student	WRN-16-2	ShuffleNet-v1	ShuffleNet-v1	ShuffleNet-v2	ResNet20	ResNet32
Teacher Acc.	76.51	76.51	79.45	79.45	73.28	74.13
Student Acc.	73.26 (0.050)	70.50 (0.360)	70.50 (0.360)	71.82 (0.062)	69.06 (0.052)	71.14 (0.061)
L2T-ww* [18]	-	-	76.35	77.39	-	-
CRD [39]	75.64	76.27	75.12	76.05	71.63	-
SSKD [47]	76.04	76.47*	77.21*	77.49*	71.49	-
ReviewKD* [4]	76.20 (0.030)	77.14 (0.015)	76.41 (0.063)	77.37 (0.069)	71.89 (0.056)	73.16 (0.029)
Equally weighted	75.50 (0.010)	74.28 (0.085)	73.54 (0.120)	74.39 (0.119)	70.89 (0.065)	73.19 (0.052)
Use α_T	76.25 (0.034)	77.19 (0.074)	77.15 (0.043)	76.64 (1.335)	71.24 (0.014)	73.58 (0.012)
DistPro	76.36 (0.005)	77.24 (0.063)	77.18 (0.047)	77.54 (0.059)	72.03 (0.022)	73.74 (0.011)

Table 1: Results on CIFAR100. Results are averaged over 5 runs. Variances are in the parentheses. “*” represents our reproduced results. In “Equally weighted”, we do not use the searched α in the retrain phase. Instead, each element of α is uniformly set to $1/L$, where L is the length of α . In “Use α_T ”, the finally converged α is used at each iteration of the retrain phase.

exactly the same training setting and hyper-parameters for all the methods, including data pre-processing, learning rate scheduler, number of training epochs, batch size, etc. We first demonstrate the effectiveness of DistPro for classification on CIFAR100. Then, we provide more analysis with a larger-scale dataset, ImageNet1K. At last, we show the results of dense prection tasks and ablation study. All experiments are performed with Tesla-V100 GPUs.

4.1 Classification on CIFAR100

Implementation details We follow the same data pre-processing approach as in [4]. Similarly, we select a group of representative network architectures including ResNet [12], WideResNet [51], MobileNet [16, 36], and ShuffleNet [52, 30]. We follow the same training setting in [4] for distillation. We train the models for 240 epochs and decay the learning rate by 0.1 for every 30 epochs after the first 150 epochs. Batch size is 128 for all the models. We train the models with the same setting five times and report the mean and variance of the accuracy on the testing set, to demonstrate the improvements are significant. As mentioned, before distillation training, we need to obtain the distillation process by searching described as below.

To search the process \mathcal{A} , we randomly split the original training set into two subdatasets, 80% of the images for training and 20% for validating. We run the search for 40 epochs and decay the learning rate for θ (parameters of the student) by 0.1 at epoch 10, 20, and 30. The learning rate for α is set to 0.05. Following the settings in [4], we did not use all feature maps for knowledge distillation but only the ones after each downsampling stage. Transform blocks are used to transform the student feature map. Fig. 4a shows the block architecture. The size of the pathways is 27, as we use 3 transform blocks and 9 connections. To make fair comparison, we follow the same HCL distillation loss in [4]. Once the the process is obtained, in order to align the searched distillation process, linear interpolation is used to expand the process of \mathcal{A} from 40 epoch (search stage) to 240 epochs (retrain stage) during KD.

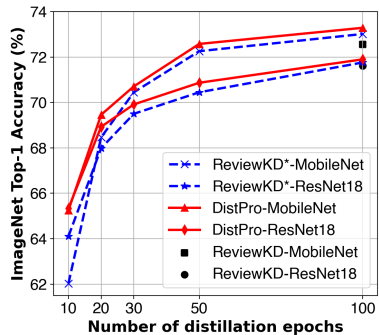


Fig. 5: Comparison study on fast distillation on ImageNet1K.

Results The quantitative results are summarized in Table 1. First two rows show the architecture of teacher and student network and their accuracy without distillation correspondingly. In the line of “ReviewKD*”, we list the results reproduced using the code released by the author [4]. To demonstrate learning the α is essential, we first assign equally weighted α to all the pathways. As shown in line “Equally weighted”, the results are worse than ReviewKD, indicating the selected pathways from ReviewKD are useful. For ablation, we first adopt the learned α_T at the end of the process. As shown in “Use α_T ”, it outperforms ReviewKD in multiple settings. The last row shows the results adopting the learned process \mathcal{A} , which outperforms ReviewKD significantly.

4.2 Classification on ImageNet1K

Implementation Details We follow the search strategy used in CIFAR100 and training configurations in [4] with a batch size of 512 (4GPUs are used). The selected student/teacher networks are MobileNet/ResNet50, ResNet18/ResNet34 and DeiT-tiny [40]/ViT-B [9]. We adopt a different architecture of transform block for DeiT. At searching stage, we search \mathcal{A} with Tiny-ImageNet [23] for 20 epochs. We adopt cosine lr scheduler for learning the student network. For search space, 1 transform block is used, while we consider 5 feature maps in the networks, and build 15 pathways by removing some pathways following [4] due to limited GPU memory. At KD stage, we train the network for various epochs with initial learning rate of 0.1, and cosine scheduler. It takes 4 GPU hours for searching and 80 GPU hours for distill with 100 epochs. More details (ViT transform block, selected pathways etc) can be found in supplementary.

Fast Distillation. From prior works [38], KD is used to accelerate the network learning especially with large dataset. This is due to the fact that its accuracy can be increased using only Ground Truth labels if the network is trained with a large number of epochs [15]. Previous works [4] only check the results stopping at 100 epochs. Here, we argue that it is also important to evaluate the results with less training cost, which is another important index for evaluating KD methods. This will be useful for the applications requiring fast learning of the network.

	Network	ReviewKD[4]	DistPro(Ours)
Top-1	MobileNet	72.56	72.54
#Epochs		100	50
Top-1	ResNet18	71.61	71.59
#Epochs		100	65
Top-1	DeiT	73.44	73.41
#Epochs		150	100

Table 2: Comparison study on numbers of distillation epochs to achieve same top-1 accuracy on ImageNet1K (lower the better).

Setting	Search dataset	Retrain dataset	Teacher	Student	Top-1 (%)
(a)	Tiny-ImageNet	ImageNet1K	ResNet50	MobileNet	73.26
	CIFAR100	ImageNet1K	ResNet50	MobileNet	73.20

Setting	Search teacher	Search student	Retrain teacher	Retrain student	Top-1 (%)
(b)	ResNet34	ResNet18	ResNet34	ResNet18	71.89
	ResNet50	MobileNet	ResNet34	ResNet18	71.87

Table 3: Results of transferring searched \mathcal{A} cross search-retrain datasets and search-retrain networks performed on ImageNet1K with 100 epochs. Top-1 accuracy on validation set is reported.

Setting	Acc. (%)	Teacher	Student	OFD[13]	LONDON[37]	SCKD[53]	ReviewKD[4]	ReviewKD*	DistPro
(a)	Top1	76.16	68.87	71.25	72.36	72.4	72.56	73.12	73.26
	Top5	93.86	88.76	90.34	91.03	-	91.00	91.22	91.27
(b)	Top1	73.31	69.75	70.81	-	71.3	71.61	71.76	71.89
	Top5	91.42	89.07	89.98	-	-	90.51	90.67	90.76
(c)	Top-1	85.1	73.41	-	-	-	-	73.44	73.51

Table 4: Comparison study on ImageNet1K. Settings are (a) teacher: ResNet-50, student: MobileNet; (b) teacher: ResNet-34, student: ResNet-18; (c) teacher: ViT-B, student: DeiT-tiny. ReviewKD* denotes our reproduced experimental results with cosine learning rate scheduler.

Here, we show **DistPro** reaches much better trade-off between training cost and accuracy. In Fig. 5, **DistPro** with different network architectures (red curves) outperforms ReviewKD with the same setting (dashed blue curves) at all proposed training epochs. The performance gain is larger when less training cost is required, *e.g.* it is 3.22% on MobileNet with 10 epochs (65.25% vs 62.03%), while decreased to 0.14% with 100 epochs, which is still a decent improvement. Similar results are observed with ResNet18. Note here for all results with various epochs, we adopt the same learning process \mathcal{A} , where the time cost can be ignored comparing to that of KD. In Tab. 2, we show the number of epochs saved by **DistPro** when achieving the same accuracy using ReviewKD [4]. For instance, in training MobileNet, **DistPro** only use 50 epochs to achieve 72.54%, which is comparable to ReviewKD trained with 100 epochs (72.56%), yielding 2x acceleration. Similar acceleration is also observed with ResNet18 and DeiT.

Transfer \mathcal{A} Here, we study whether learned \mathcal{A} can be transferred across datasets and similar architectures. Tab. 3 shows the results. In setting (a), we resize the image in CIFAR100 to 224×224 , and do the process search on CIFAR100, where the search cost is only 2 GPU hours. As shown at last column, it only downgrades accuracy by 0.06%, which is comparable with full results searched with ImageNet, demonstrating the process could be transferred across datasets. In setting (b), we adopt the searched process with student/teacher networks of MobileNet/ResNet50 to networks of ResNet18/ResNet34 since they have same feature pathways at similar corresponding layers. As shown, the results are closed, demonstrating the process could be transferred.

Best Results. Finally, Tab.4 shows the quantitative results comparing with various SoTA baselines. As mentioned, we adopt cosine scheduler while ReviewKD [4] adopts step scheduler. To make it fair, we retrain ReviewKD with cosine scheduler, and list the results in the column of ReviewKD*. As shown in

Teacher	ResNet101
Student	MobileNet-v2
Baseline	66.92 (0.00721)
IFVD [44]	68.31 (0.00264)
SCKD [53]	68.25 (0.00307)
+ReviewKD [4]	69.03 (0.00373)
+Equally Weighted α	68.49 (0.00793)
+DistPro	69.12 (0.00462)

Table 5: mIoU (%) on CityScapes (higher the better). Results are averaged over 5 runs. Standard deviation in the parentheses.

Teacher	ResNet50
Student	ResNet18
Baseline	0.2032
KD	0.2045
ReviewKD [4]	0.1983
Equally Weighted α	0.2030
DistPro	0.1972

Table 6: Estimation error on NYUv2 (lower the better).

Normalization	Acc.
$\text{softmax}(\text{concat}(\tilde{\alpha}, 1))$	79.78
$\text{softmax}(\tilde{\alpha})$	79.41
$\frac{\tilde{\alpha}}{\ \tilde{\alpha}\ _1 + 1}$	79.47
$\frac{\tilde{\alpha}}{\ \tilde{\alpha}\ _1}$	79.20
$\text{sigmoid}(\tilde{\alpha})$	78.80

Table 7: Performance of different normalization methods on CIFRA-100. $\text{concat}(\tilde{\alpha}, 1)$ concatenate $\tilde{\alpha}$ with a scalar 1. Results are averaged over 5 runs.

the table, for three distillation settings, DistPro outperform the existing methods achieving top-1 accuracy of 73.26% for MobileNet, 71.89% for ResNet18 and 73.51% for DeiT respectively, yielding new SoTA results for all these networks.

4.3 Dense prediction tasks

CityScapes is a popular semantic segmentation dataset with pixel class labels [6]. We compare to a SoTA segmentation KD method IFVD [44] and adopt their released code. IFVD is a response-based KD method, therefore can be combined with feature-based distillation method including ReviewKD and DistPro, which are shown as “+ReviewKD” and “+DistPro” respectively in Table 5. We adopt student of MobileNet-v2 to compare against another SoTA results from SCKD [53]. Please note here for simplicity and numerical stability, we disable adversarial loss of the original IFVD. From the table, “+DistPro” outperforms all competing methods.

NYUv2 is a dataset widely used for depth estimation [32]. The experiments are based on the code released in S2D [29]. We compare DistPro with the plain KD, where we treat teacher output as ground truth, and ReviewKD with intermediate feature maps. Root mean squared errors (RMSE) are summarized in Table 6, and it shows that DistPro is also beneficial.

4.4 Ablation study

Is it always beneficial to transfer knowledge only from lower-level feature maps to higher-level feature maps? In Table 6 of [4], the authors conducted a group of experiments on CIFAR100, which show that only pathways from lower-level feature maps in teacher to higher-level feature maps in student are beneficial. However, we conducted similar experiments on CityScapes, the results did not support this claim. Specifically, using pathways from the last feature map to the first three feature maps results in mIoU 73.9% on ResNet18, while using pathways from all lower-level feature maps to higher level ones as in [4] results in 73.5%. This suggests that the optimal teaching scheme need to be searched w.r.t larger space. Also, the results in Table 5 also show that our searched process is better than the hand-crafted one.

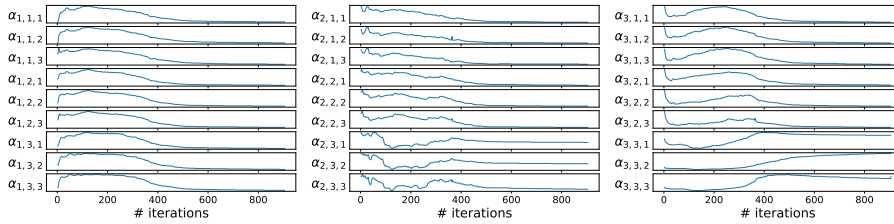


Fig. 6: Searched distillation process performed on CIFAR100 with WRN-40-2 as teacher and WRN-16-2 as student. x -axis is the number of iterations. The left image contains $\alpha[1, :, :]$, i.e., all the elements corresponding to the lowest-level feature map of the teacher. The middle one corresponds to $\alpha[2, :, :]$, and the right one corresponds to $\alpha[3, :, :]$.

Intuition of process \mathcal{A} In Figure 6, we show how the learned sample of \mathcal{A} changes with time in search stage of distilling WRN-40-2 to WRN-16-2. Similar observations are found in other settings. The figure indicates that at the early stage of the training, the optimized teaching scheme focuses on transferring knowledge from low-level feature maps of the teacher to the student. As training goes, the optimized teaching scheme gradually moves on to the higher-level feature maps of the teacher. Intuitively, high-level feature maps encode highly abstracted information of the input image and thus are harder to learn from compared to the low-level feature maps. **DistPro** is able to automatically find a teaching scheme to make use of this intuition.

Normalization of α_t As mentioned before, we apply normalization to α_t ((2)) for numerical stability after 2_{nd} gradient approximation. In this study, we evaluate several normalization strategies. Let the unnormalized parameters be $\tilde{\alpha}_t$ and the normalized ones be α_t . In the following, we view α_t as a vector of length L . The experiments are conducted on CIFAR100 with WRN-28-4 as teacher and WRN-16-4 as student. The results are shown in Table 7, and the proposed biased softmax normalization outperforms others. Intuitively, due to the appended scalar 1, the value of $\tilde{\alpha}_t$ can be compared against 1, and a constant value yields a label smoothing [31] effect for the distribution.

5 Conclusion

In this paper, we take a step towards the problem of finding the optimal KD scheme given a pair of wanted student network and learned teacher network under a vision task. Specifically, we setup a searching space by building pathways between the two networks and assigning a stochastic distillation process along each pathway. We propose a meta-learning framework, **DistPro**, to learn these processes efficiently, and find effective ones to perform KD with intermediate features. We demonstrate its benefits over image classification, and dense predictions such as image segmentation and depth estimation. We hope our method could inspire the field of KD to further expand the scope, and its cooperation with other techniques such as NAS and hyper-parameter tuning.

References

1. Andrychowicz, M., Denil, M., Gomez, S., Hoffman, M.W., Pfau, D., Schaul, T., Shillingford, B., De Freitas, N.: Learning to learn by gradient descent by gradient descent. In: *Advances in neural information processing systems*. pp. 3981–3989 (2016)
2. Buciluă, C., Caruana, R., Niculescu-Mizil, A.: Model compression. In: *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. pp. 535–541 (2006)
3. Chen, G., Choi, W., Yu, X., Han, T., Chandraker, M.: Learning efficient object detection models with knowledge distillation. *Advances in neural information processing systems* **30** (2017)
4. Chen, P., Liu, S., Zhao, H., Jia, J.: Distilling knowledge via knowledge review. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 5008–5017 (2021)
5. Chen, Z., Liu, B.: Lifelong machine learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning* **12**(3), 1–207 (2018)
6. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016)
7. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A Large-Scale Hierarchical Image Database. In: *CVPR09* (2009)
8. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018)
9. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020)
10. Franceschi, L., Frasconi, P., Salzo, S., Grazzi, R., Pontil, M.: Bilevel programming for hyperparameter optimization and meta-learning. In: *International Conference on Machine Learning*. pp. 1568–1577. PMLR (2018)
11. Gou, J., Yu, B., Maybank, S.J., Tao, D.: Knowledge distillation: A survey. *International Journal of Computer Vision* **129**(6), 1789–1819 (2021)
12. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770–778 (2016)
13. Heo, B., Kim, J., Yun, S., Park, H., Kwak, N., Choi, J.Y.: A comprehensive overhaul of feature distillation. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 1921–1930 (2019)
14. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network (2015)
15. Hoffer, E., Hubara, I., Soudry, D.: Train longer, generalize better: closing the generalization gap in large batch training of neural networks. *Advances in neural information processing systems* **30** (2017)
16. Howard, A.G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., Adam, H.: Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861* (2017)

17. Huang, Z., Wang, N.: Like what you like: Knowledge distill via neuron selectivity transfer. arXiv preprint arXiv:1707.01219 (2017)
18. Jang, Y., Lee, H., Hwang, S.J., Shin, J.: Learning what and where to transfer. In: International Conference on Machine Learning. pp. 3030–3039. PMLR (2019)
19. Ji, M., Heo, B., Park, S.: Show, attend and distill: Knowledge distillation via attention-based feature matching. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 7945–7952 (2021)
20. Kimura, A., Ghahramani, Z., Takeuchi, K., Iwata, T., Ueda, N.: Few-shot learning of neural networks from scratch by pseudo example optimization. arXiv preprint arXiv:1802.03039 (2018)
21. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
22. Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images (2009)
23. Le, Y., Yang, X.: Tiny imagenet visual recognition challenge. CS 231N **7**(7), 3 (2015)
24. Lee, S., Song, B.C.: Graph-based knowledge distillation by multi-head attention network. arXiv preprint arXiv:1907.02226 (2019)
25. Liu, H., Simonyan, K., Yang, Y.: Darts: Differentiable architecture search. In: International Conference on Learning Representations (2019)
26. Liu, Y., Chen, K., Liu, C., Qin, Z., Luo, Z., Wang, J.: Structured knowledge distillation for semantic segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2604–2613 (2019)
27. Liu, Y., Jia, X., Tan, M., Vemulapalli, R., Zhu, Y., Green, B., Wang, X.: Search to distill: Pearls are everywhere but not the eyes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7539–7548 (2020)
28. Lyu, L., Chen, C.H.: Differentially private knowledge distillation for mobile analytics. In: Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval. pp. 1809–1812 (2020)
29. Ma, F., Karaman, S.: Sparse-to-dense: Depth prediction from sparse depth samples and a single image. In: 2018 IEEE international conference on robotics and automation (ICRA). pp. 4796–4803. IEEE (2018)
30. Ma, N., Zhang, X., Zheng, H.T., Sun, J.: Shufflenet v2: Practical guidelines for efficient cnn architecture design. In: Proceedings of the European conference on computer vision (ECCV). pp. 116–131 (2018)
31. Müller, R., Kornblith, S., Hinton, G.E.: When does label smoothing help? Advances in neural information processing systems **32** (2019)
32. Nathan Silberman, Derek Hoiem, P.K., Fergus, R.: Indoor segmentation and support inference from rgb-d images. In: ECCV (2012)
33. Oord, A., Li, Y., Babuschkin, I., Simonyan, K., Vinyals, O., Kavukcuoglu, K., Driessche, G., Lockhart, E., Cobo, L., Stimberg, F., et al.: Parallel wavenet: Fast high-fidelity speech synthesis. In: International conference on machine learning. pp. 3918–3926. PMLR (2018)
34. Ren, M., Zeng, W., Yang, B., Urtasun, R.: Learning to reweight examples for robust deep learning. In: International Conference on Machine Learning. pp. 4334–4343. PMLR (2018)
35. Romero, A., Ballas, N., Kahou, S., Chassang, A., Gatta, C., Bengio, Y.: Fitnets: Hints for thin deep nets. CoRR **abs/1412.6550** (2015)
36. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C.: Mobilenetv2: Inverted residuals and linear bottlenecks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4510–4520 (2018)

37. Shang, Y., Duan, B., Zong, Z., Nie, L., Yan, Y.: Lipschitz continuity guided knowledge distillation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10675–10684 (2021)
38. Shen, Z., Xing, E.: A fast knowledge distillation framework for visual recognition. arXiv preprint arXiv:2112.01528 (2021)
39. Tian, Y., Krishnan, D., Isola, P.: Contrastive representation distillation. In: International Conference on Learning Representations (2020)
40. Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jégou, H.: Training data-efficient image transformers & distillation through attention. In: International Conference on Machine Learning. pp. 10347–10357. PMLR (2021)
41. Tung, F., Mori, G.: Similarity-preserving knowledge distillation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 1365–1374 (2019)
42. Urner, R., Shalev-Shwartz, S., Ben-David, S.: Access to unlabeled data can speed up prediction time. In: ICML (2011)
43. Wang, X., Zhang, R., Sun, Y., Qi, J.: Kdgan: Knowledge distillation with generative adversarial networks. In: NeurIPS. pp. 783–794 (2018)
44. Wang, Y., Zhou, W., Jiang, T., Bai, X., Xu, Y.: Intra-class feature variation distillation for semantic segmentation. In: European Conference on Computer Vision. pp. 346–362. Springer (2020)
45. Xiang, L., Ding, G., Han, J.: Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification. In: European Conference on Computer Vision. pp. 247–263. Springer (2020)
46. Xie, S., Zheng, H., Liu, C., Lin, L.: SNAS: stochastic neural architecture search. In: International Conference on Learning Representations (2019), <https://openreview.net/forum?id=rylqooRqK7>
47. Xu, G., Liu, Z., Li, X., Loy, C.C.: Knowledge distillation meets self-supervision. In: European Conference on Computer Vision (ECCV) (2020)
48. Yao, L., Pi, R., Xu, H., Zhang, W., Li, Z., Zhang, T.: Joint-detnas: Upgrade your detector with nas, pruning and dynamic distillation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10175–10184 (2021)
49. Yim, J., Joo, D., Bae, J., Kim, J.: A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4133–4141 (2017)
50. Zagoruyko, S., Komodakis, N.: Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. arXiv preprint arXiv:1612.03928 (2016)
51. Zagoruyko, S., Komodakis, N.: Wide residual networks. arXiv preprint arXiv:1605.07146 (2016)
52. Zhang, X., Zhou, X., Lin, M., Sun, J.: Shufflenet: An extremely efficient convolutional neural network for mobile devices. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 6848–6856 (2018)
53. Zhu, Y., Wang, Y.: Student customized knowledge distillation: Bridging the gap between student and teacher. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5057–5066 (2021)