

This paper is a preprint (Accepted in ICPR2022).

©2022 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

arXiv:2209.07606v1 [cs.CV] 15 Sep 2022

# CES-KD: Curriculum-based Expert Selection for Guided Knowledge Distillation

Ibtihel Amara  
McGill University  
Montreal, Canada

ibtihel.amara@mail.mcgill.ca

Maryam Ziaefard  
McGill University  
Montreal, Canada

maryam.ziaefard@mcgill.ca

Brett H. Meyer  
McGill University  
Montreal, Canada

brett.meyer@mcgill.ca

Warren Gross  
McGill University  
Montreal, Canada

warren.gross@mcgill.ca

James J. Clark  
McGill University  
Montreal, Canada

james.j.clark@mcgill.ca

**Abstract**—Knowledge distillation (KD) is an effective tool for compressing deep classification models for edge devices. However, the performance of KD is affected by the large capacity gap between the teacher and student networks. Recent methods have resorted to a multiple teacher assistant (TA) setting for KD, which sequentially decreases the size of the teacher model to relatively bridge the size gap between these models. This paper proposes a new technique called *Curriculum Expert Selection for Knowledge Distillation (CES-KD)* to efficiently enhance the learning of a compact student under the *capacity gap problem*. This technique is built upon the hypothesis that a student network should be guided gradually using stratified teaching curriculum as it learns easy (hard) data samples better and faster from a lower (higher) capacity teacher network. Specifically, our method is a gradual TA-based KD technique that selects a single teacher per input image based on a curriculum driven by the difficulty in classifying the image. In this work, we empirically verify our hypothesis and rigorously experiment with CIFAR-10, CIFAR-100, CINIC-10, and ImageNet datasets and show improved accuracy on VGG-like models, ResNets, and WideResNets architectures.

## I. INTRODUCTION

Modern deep networks are over-parameterized and are computationally expensive to be deployed onto edge devices. There have been a tremendous focus on compressing large models in the literature [1]–[6]. Knowledge distillation (KD) [7], a model compression technique, which relies on a teacher-student training protocol, has gained popularity in the past few years. The success of KD is mainly associated to its noticeable versatility and generalization aspect. In other words there are no constraints on the type of network architecture. Instead, “any teacher model can teach any student” [8], to a certain extent.

However, “*every success hides in it some multiple shortcomings*”. In point of fact, it has been shown in [9] that performing KD does not always yield better student performance. KD might not succeed if the capacity of the student network is much lower than the teacher’s capacity. This phenomenon is called the *capacity gap problem* in KD. Multiple teacher assistants (TA) of intermediate capacity sizes has become a go-to technique to overcome this shortcoming [9], [10]. However the sequential distillation process in [9] and the densely guided ensemble distillation process for learning TAs [10] might not help a lot in enhancing the compressed student’s performance since the former depends solely on a single TA network and the latter exploits the average (i.e. aggregated) knowledge of the

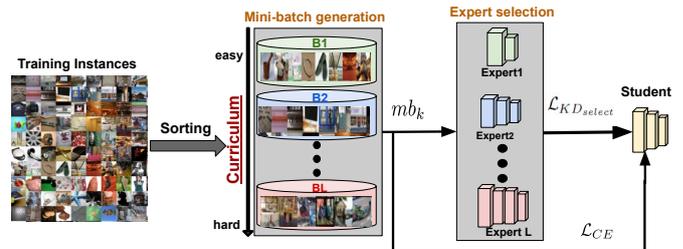


Fig. 1. Overall Pipeline of the proposed CES-KD framework. Given a training dataset, we design an easy-to-hard curriculum on the data using a meta-network (described in III-A). Then, we bucketize the sorted training set according to the total number of available teacher and assistant networks (i.e. experts). The student network is guided via distillation by selecting a single expert within the pool of teacher and assistant models. This expert selection is determined by  $\mathcal{L}_{KD_{select}}$  loss between the expert and the student. The total loss of CES-KD pipeline consists of both the distillation loss  $\mathcal{L}_{KD_{select}}$  and the cross entropy loss between the student and the sorted data  $\mathcal{L}_{CE}$ .

TAs and ignores the diversity and importance of each expert (i.e. TA networks) within the ensemble.

Intuitively, the learning process of a compressed student network can benefit from a stratified teaching and sub-curriculum oriented KD. Indeed, the student’s training behaviour changes when faced with easy-to-hard curriculum on the data [11] and also when faced with different teacher or experts in terms of size for distillation. Therefore, we hypothesize that a compact student network has faster and better learning ability on easy data samples when guided by a less complex or lower capacity teacher network. Similarly, a compact student learns better and faster on difficult data samples when guided by a more complex or higher capacity teacher model during training.

We propose *Curriculum Expert Selection for Knowledge Distillation (CES-KD)*, a TA-based KD technique that is established on a single expert selection per input sample to guide the KD process from a large cumbersome teacher network down to a compressed low capacity student network. Our method borrows insights from the field of curriculum learning and adopts the multiple teaching assistant scheme [9], [10], but mainly leverages the individuality of each TA network within the ensemble. In particular, our curriculum learning approach can be globally summarized as follows: given the level of classification difficulty of an input image, we assign an

appropriate expert (i.e teacher network) for distillation. More details can be found in Section III. The main contributions of our paper are:

- 1) We propose a curriculum based KD approach that intuitively guides the learning process of a compact student network by selecting the appropriate teacher assistant network according to the provided input samples.
- 2) We empirically show that on easy data samples, the compressed student learns better and faster from low capacity teachers and on difficult data samples the student learns better and faster from higher capacity teachers.
- 3) We improve the accuracy of the student network on various datasets and architectures as compared to baseline and state of the art methods.
- 4) We show that our method has faster convergence than state of the art methods.

## II. RELATED WORKS

### A. Knowledge Distillation

Knowledge distillation is a training method that is based on teacher-student learning. It was firstly introduced as a mode of model compression by Bucilua et al. [12] then further popularized by Hinton et al. [7]. The goal of KD is to increase the accuracy of a student network due to transfer of information from the pre-trained teacher network during training. Distillation methods rely on different techniques to capture the knowledge of the teacher that is transferred to the student. The traditional KD [7] uses the soft targets produced by the teacher network as the knowledge to transfer to the student. Some recent works focus on feature-based knowledge distillation and train the student to match the intermediate layers of the teacher [13]. In addition, some approaches [14] transfer spatial attention maps, where the student attends to similar parts of the image as the teacher network. In our work we exploit the soft targets of the teacher network for our distillation process.

### B. Capacity gap problem

It was empirically shown in [9] that compressed student networks are harder to train when the size gap between the teacher and the student model is very large. Many papers have attempted to solve this issue [9], [10], [15]. Notably KD techniques that used multiple teacher assistant (TA) networks have stood out to relatively solve the capacity gap problem in KD. TAKD [9] proposed using intermediate capacity networks to distill Knowledge from a large dense network down to a compact student. DGKD [10] is a method to mitigate the capacity gap problem through densely performing the TAKD process for knowledge distillation. In this way each teacher assistant will learn from the ensemble of previous teacher assistants and teachers. Similarly, the student will learn from all predefined teacher assistants and the main teacher. Our method differs from these techniques as we use a curriculum-based paradigm to train the student network, using the generated multiple assistant models from the large teacher network. We

explicitly utilize the obtained teaching assistant models and the large teacher model to guide the learning of the student. Unlike the work in [10], we do not aggregate the knowledge of these multiple assistants and teacher to train the student model. Instead, we exploit the diversity of these networks and their individual expertise on the training data through the expert selection aspect of our method.

### C. Curriculum learning in KD

The field of curriculum learning (CL) became popular in deep learning due to its ability to alleviate certain training problems by tackling the structure of the training data instead of the network architecture. CL imposes an order on the training process of a student network. It was shown that this technique efficiently enhances the performance of deep models [16]. The CL process consists of two major steps: (1) a scoring function, which organizes the data by level of difficulty and (2) a pacing function, which defines the process of feeding the sorted data into the network. There are various methods to score and sort a dataset and to perform pacing functions. These were investigated by Hacothen and Weinshall [17]. CL was also shown to be beneficial for KD. Xiang et al. [18] trained a student network and implemented a self-paced learning function to classify imbalanced datasets and showed a substantial boost in performance. Zhao et al. [19] performs an instance-level KD where the curriculum is applied on the data and is ensured by a snapshot copy of a student network. Panagiotatos et al. [20] applied curriculum on the teacher network rather than on the data. They showed that teachers of different learning levels can guide the compact student's training during distillation. In their work, they took different versions of a teacher network at different training point and used them to perform ensemble knowledge distillation on the entire dataset. In our work, we focus on applying curriculum on both the dataset and the teacher models. We sort the dataset based on an easy-to-hard curriculum. Then, we perform a bucketing mechanism, which divides the sorted dataset into different buckets according to the available teacher and teacher assistants. Each bucket is assigned to a specialized teacher assistant or teacher network based on its level of expertise. More details can be found in Section III.

## III. CURRICULUM-BASED EXPERT SELECTION FOR KD

### A. Hypothesis and Motivation

We hypothesize that a student network learns better and faster from a small TA network when faced with easy concepts and better and faster from a large teacher network when dealing with difficult concepts. To validate this hypothesis, we trained a compact student network (ResNet20) on CINIC-10 dataset. The training data are sorted according to an easy-to-hard curriculum using a finetuned meta-network (full details on designing this curriculum is given in Section III-B1. We then divide the obtained curriculum into three separate subsets of different levels of difficulty: Easy, Intermediate and Difficult. We also took different teacher networks of different sizes (i.e capacities) to guide the student's training via baseline

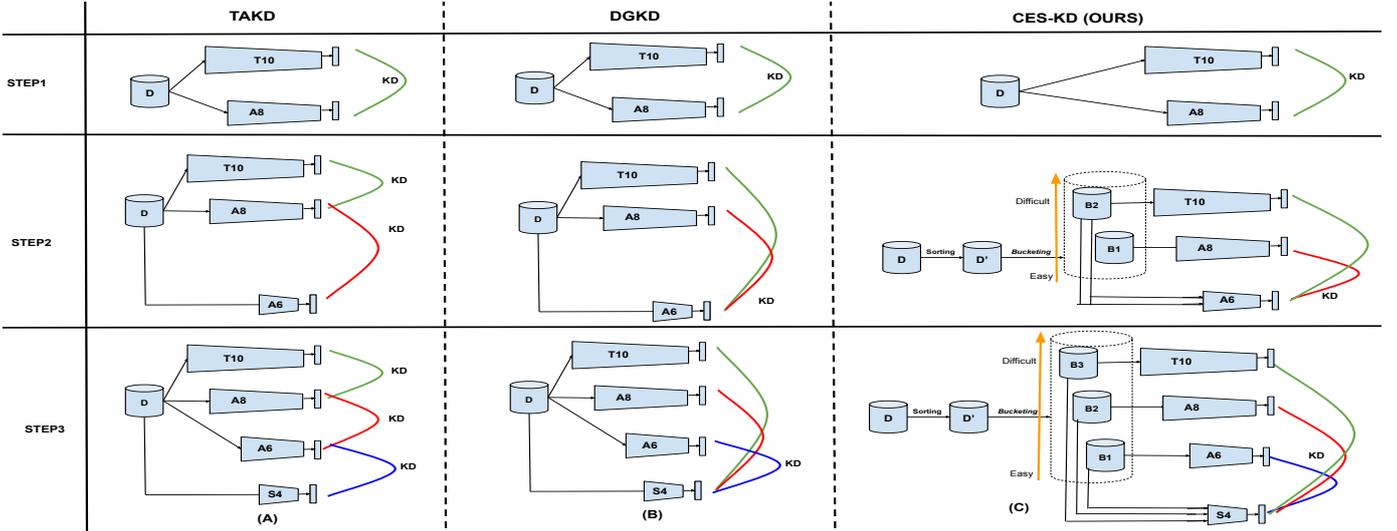


Fig. 2. **Back-to-back comparison of the distillation process of TAKD (first column), DGKD (middle column) and CES-KD (last column) methods.** Three-step distillation process from a teacher network of 10 layers (T10) down to a student network of 4 layers (T4), with intermediate teacher assistants A8 and A6. All three methods share the first distillation step giving the assistant network A8. As for the second distillation, TAKD performs a baseline KD from A8 to A6. DGKD performs ensemble KD of T10 and A8 to A6. CES-KD (our method) is a two-step process. First, we sort the data samples by level of difficulty using a scoring function (This process is fixed and is done only once during the entire distillation processes). Second, we bucket the sorted training dataset (equal division of samples into buckets) and assign each bucket to a designated expert (i.e. teacher / assistant network). The bucket containing the easiest data samples is given to the latest teacher assistant and the bucket containing the most difficult samples in the dataset are set to the large teacher network. This bucketing technique is dynamic: at each distillation step the sorted data is equally divided to the total number of experts. (For example two assistants and one teacher network yields to three buckets in total.)

knowledge distillation (BLKD) [7] on these different levels of curriculum. Table I shows the test accuracy of the student network trained on three different levels of difficulty with three teachers of different capacities. We see that for easy samples the student network (ResNet20) has higher accuracy when guided through distillation by the lowest capacity teacher (ResNet26) within the group of experts. As for the difficult samples, the student acquires better knowledge from the highest capacity teacher model (ResNet56), which validates our assumption in terms of quality of learning. To further study the student’s learning efficiency, Figure 3 (a) shows that the student’s optimization, on easy samples, converges faster when trained with the guidance of the lowest capacity teacher model (ResNet26) via knowledge distillation. Also in Figure 3 (b) we observe that on difficult samples, the student learns faster (i.e. faster convergence) from the highest capacity network (ResNet56) than from the lower capacity teachers. This comparison validates our hypothesis and motivates our technique regarding the curriculum data-model selection.

TABLE I  
TOP-1 % TEST ACCURACY OF THE STUDENT NETWORK (RESNET20) TRAINED ON THREE CURRICULUM LEVELS UNDER THE SUPERVISION (DISTILLATION) OF DIFFERENT CAPACITY TEACHER NETWORKS ON THE TEST SET OF CINIC10 DATASET. AVERAGE OVER THREE INDEPENDENT RUNS.

Teacher networks	Easy	Intermediate	Difficult
Resnet 56	79.08 %	81.59 %	<b>79.20%</b>
Resnet 32	79.28 %	<b>81.88 %</b>	79.10 %
Resnet 26	<b>79.30 %</b>	<b>81.88 %</b>	78.97 %

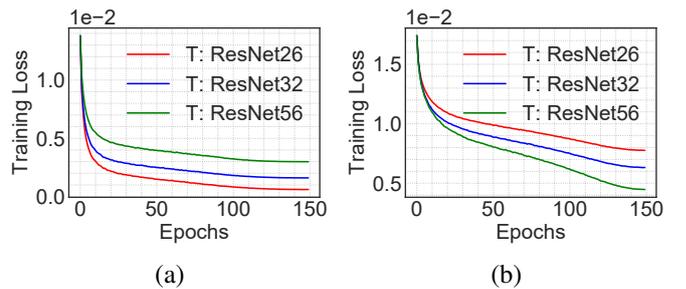


Fig. 3. Training loss of a compact student network (S:ResNet20) distilled using baseline knowledge distillation technique (BLKD) on easy data samples (a) and on hard data samples (b) from different teacher networks of different capacities (T: ResNet26, ResNet32, ResNet56).

## B. Methodology of our proposed method

Figure 1 presents a global overview of the distillation framework CES-KD. Our method relies on two main steps: the design of the data curriculum and the selection of a single representative expert based on their expertise on a given input data (i.e. images). We provide in Figure 2 details on the distillation pipeline of our method with a back-to-back comparison on current TA-based KD methods (TAKD [9] and DGKD [10]). Following the standard CL paradigm, we need to address two main questions: (1) How do we rank the dataset according to an easy-to-hard curriculum? (see subsection III-B1); (2) How do we train a student model using the ranked instances? (See subsection III-B2).

1) *Design of the data curriculum:* In this work, we adopt the *meta-network* method described by Hacohen and Weinschall [17] as *transfer learning-based scoring function*. In

particular, we consider a reference model trained on a very large dataset. Then, we fine-tune this reference model on the smaller training dataset. We evaluate the difficulty score of the training data using the real-valued loss given by the fine-tuned meta-network. This method was also previously explored in [21]. Therefore, for a given reference model with weights  $\mathbf{W}$ ,  $f_{\mathbf{W}} : \mathbf{X} \rightarrow \mathbf{Y}$ , the difficulty score of a sample  $x_i$  is defined by  $s(x_i, y_i) = \text{Loss}(f_{\mathbf{W}}(x_i), y_i)$ . Finally, these samples are sorted in an ascending order based on their score values, resulting in a sorted dataset called  $D_{\text{sorted}}$ . This curriculum specification is performed only once during our entire distillation pipeline.

2) *Representative expert selection*: After ranking the instances within the training dataset, we perform a *bucketing system* in which we split  $D_{\text{sorted}}$  into  $L$  equal buckets  $\{B_1, \dots, B_L\}$  according to the number of experts (teacher and teacher assistants) we are considering for distillation, at a particular distillation step. The first bucket contains the easiest samples and the last bucket includes the hardest samples. Each bucket is assigned to a representative teacher or TA network based on their level of expertise. The first bucket containing the easiest samples are given to the lowest capacity teacher network. The next bucket is given to the second lowest capacity teacher assistant network, until the last bucket containing the hardest data samples are sent to the largest teacher network.

---

#### Algorithm 1 CES-KD

---

- 1: Rank the original dataset  $D$  according to scoring function  $s: D_{\text{ranked}}$
  - 2: Divide  $D_{\text{ranked}}$  into  $L$  equal buckets:  $\{B_1, \dots, B_L\}$ .
  - 3: Assign the experts in an ascending order of depths to the buckets, i.e the shallowest expert is assigned to  $B_1$  until the deepest expert is assigned to  $B_L$
  - 4: Generate mini-batches  $mb_k$  in each bucket
  - 5: **for** all mini-batches in  $D_{\text{ranked}}$  **do**
  - 6:     **for** each bucket  $B_l$  in  $\{B_1, \dots, B_L\}$  **do**
  - 7:         **if**  $mb_k$  belongs to  $B_l$  **then**
  - 8:             Select the representative expert  $l$  from the ensemble to provide soft targets for  $mb_k$
  - 9:             **end if**
  - 10:         **end for**
  - 11:     Update student weights with  $mb_k$
  - 12: **end for**
- 

3) *Distillation loss*: Following the proposed data curriculum and teacher selection, we train the student network according to the ascending level of difficulty of instances in  $D_{\text{ranked}}$ . We generate mini-batches  $mb_k$  within each bucket  $B_l$ , where  $k$  is the number of mini-batches per bucket, and guide the student with the output logits produced by the selected expert.  $l$  is the index of the representative teacher or TA network. We define the distillation loss as:

$$\mathcal{L}_{KD_{\text{select}}} = \sum_{l=1}^L \omega_l^{mb_k \in B_l} \mathcal{L}_{CE}(\sigma(\frac{\mathbf{z}_s}{T}), \sigma(\frac{\mathbf{z}_l}{T})) \quad (1)$$

where:

$$\omega_l^{mb_k \in B_l} = \mathbf{1}_{B_l}(mb_k) = \begin{cases} 1 & \text{if } mb_k \in B_l \\ 0 & \text{else} \end{cases}$$

$\mathcal{L}_{CE}(\cdot, \cdot)$  denotes the cross entropy loss.  $\sigma(\cdot)$  is the softmax function.  $\mathbf{z}_s$  and  $\mathbf{z}_l$  are output logits of the student and selected expert  $l$ , respectively.  $T$  is the temperature hyperparameter.  $\omega_l^{mb_k \in B_l}$  denotes a weight that selects the representative experts from the ensemble according to the input mini-batch.

4) *Training*: The total training loss is then defined as:

$$\mathcal{L}_{\text{CES-KD}} = \alpha T^2 \mathcal{L}_{KD} + (1 - \alpha) \mathcal{L}_{CE}(\hat{\mathbf{y}}, \sigma(\mathbf{z}_s)) \quad (2)$$

$\hat{\mathbf{y}}$  is the one-hot vector indicating the ground-truth class.  $\mathbf{z}_s$  is the output logit vector of the student model. The pseudo code of our overall methodology is provided in Algorithm 1.

## IV. EXPERIMENTAL SET-UP

### A. Datasets and Networks

We perform experiments on different datasets - CIFAR-10 [22], CINIC-10 [23], CIFAR-100 [22], and ImageNet [?], having different numbers of classes ranging from 10, 100 up to 1000. We also used different network architectures to validate the performance of our method from plain convolutions (plain CNNs), which are VGG-like networks [24], ResNets [25], and WideResNets [26]. We report the performance in terms of accuracy on the validation set for ImageNet and on the test set for the remaining of the datasets.

### B. Implementation Details

All implementations are done using PyTorch [27]. For all experiments we perform data augmentation, specifically random crops and random horizontal flips. Then we perform data normalization by subtracting the mean and the variance of the entire training set. For the experiments on plain CNN architectures on CIFAR-100 and ResNets on CIFAR-10, we used stochastic gradient descent (SGD) optimizer with Nesterov momentum of 0.9, weight decay of 1e-4, and with a mini-batch size of 128. The initial learning rate was 0.1, then divided by 10 at the 30th, 90th, and 120th epochs; we trained for a total of 150 epochs. Unlike TAKD [9] and DGKD [10], we do not use a hyper-parameter optimization toolkit on a hyper-parameter search space and seed setting. Instead, we perform multiple runs of the experiment and report the average across multiple random seeds, and the corresponding standard deviations. For fairer comparison between techniques we used the same hyper-parameters for each distillation step and across different TA-based KD techniques. As for KD related hyper-parameters, we took  $\alpha$  to be 0.9 and the temperature  $T$  to be 10. For the experiments related to benchmarking on CIFAR-100 and comparing our method to state-of-the-art KD techniques, we use the same set-up as Tian et al. in [28] detailed in their code<sup>1</sup>. For the ImageNet experiments, we use the same training set-up as the Imagenet distributed training from Pytorch<sup>2</sup>.

<sup>1</sup><https://github.com/HobbitLong/RepDistiller>

<sup>2</sup><https://github.com/pytorch/examples/tree/master/imagenet>

### C. Data Sorting and Curriculum

For experiments on CIFAR-10, CINIC-10, and CIFAR-100, we use a reference model of ResNet110 pre-trained on ImageNet. The ResNet model was taken from the pretrained torch model zoo in Pytorch [27]. We freeze the feature layers and fine-tune using only the classification layers on the targeted smaller training dataset. For the experiments related to ImageNet, we do not perform fine-tuning. Instead we consider the pre-trained network itself from [27] as the reference model. We ran extensive trials using the ranked data instances to train the student. We report the results of the curriculum that works best in training our students. All buckets are incrementally presented per epoch to the student network to avoid the issue of catastrophic forgetting [29], [30]. The classes are distributed almost equally across the buckets to avoid any biases during training. We generate mini-batches per bucket with uniform sampling of instances within the bucket.

## V. RESULTS AND DISCUSSION

### A. Comparing CES-KD to current TA-based KD methods

In this section we show the effectiveness of our proposed method CES-KD on several standard datasets such as CIFAR-10 and CIFAR-100 when compared to TA-based KD methods.

TABLE II

TEST ACCURACY WITH ALL DISTILLATION STEPS USING PLAIN CNN ARCHITECTURE ON CIFAR100 OVER THREE RANDOM SEEDS. WE ALSO REPORT THE CORRESPONDING STANDARD DEVIATION ON THE STEPS CONTAINING TA NETWORKS. TEACHER  $T_{10}$  ASSISTANTS  $A_8, A_6$  AND STUDENT  $S_4$ . (\*) METHODS USE PUBLICLY PROVIDED CODE.

Step	TAKD*	DGKD*	CES-KD
Teacher (CNN-10; $T_{10}$ )		66.89	
Student (CNN-4; $S_4$ )		62.71	
$T_{10} \rightarrow A_8$		64.20	
$T_{10} \rightarrow A_8 \rightarrow A_6$	67.80 $\pm$ 0.35	68.56 $\pm$ 0.282	<b>68.90 <math>\pm</math> 0.196</b>
$T_{10} \rightarrow A_8 \rightarrow A_6 \rightarrow S_4$	63.31 $\pm$ 0.145	63.44 $\pm$ 0.127	<b>63.58 <math>\pm</math> 0.045</b>

Table II shows the test accuracy of a student network (4-layered CNN) when distilled from a teacher network (10-layered CNN) at each level of a defined distillation path. To distill knowledge from a teacher  $T_{10}$  down to a student  $S_4$ , we performed the following distillation path  $T_{10} \rightarrow A_8 \rightarrow A_6 \rightarrow S_4$ . Our method shows good improvements overall. CES-KD achieves an accuracy of 68.90 % for the  $T_{10} \rightarrow A_8 \rightarrow A_6$  path, which is almost 1% improvement to TAKD, a notable increase compared to DGKD, and especially a 6.19 % improvement compared to the student network trained individually and from scratch. As for the path  $T_{10} \rightarrow A_8 \rightarrow A_6 \rightarrow S_4$ , we also see a substantial accuracy increase with our proposed method.

Table III shows the test accuracy of a compact student network using a residual architecture. The teacher network is a ResNet26 and the student is a ResNet8. Similarly to our previous observations, our method shows a considerable improvement when compared to both TAKD and DGKD. For both paths  $T_{26} \rightarrow A_{20} \rightarrow A_{14}$ ,  $T_{26} \rightarrow A_{20} \rightarrow A_{14} \rightarrow A_8$ , CES-KD shows, respectively, a 5.46% and a 1.62% improvement over the student network trained from scratch.

TABLE III

TEST ACCURACY WITH ALL DISTILLATION STEPS USING RESNET ARCHITECTURE ON CIFAR10 OVER THREE RANDOM SEED. WE ALSO REPORT THE CORRESPONDING STANDARD DEVIATION ON THE STEPS CONTAINING TA NETWORKS. TEACHER  $T_{26}$  ASSISTANTS  $A_{20}, A_{14}$  AND STUDENT  $S_8$ . (\*) METHODS USE PUBLICLY PROVIDED CODE.

Step	TAKD*	DGKD*	CES-KD
Teacher (ResNet26; $T_{26}$ )		91.73	
Student (ResNet8; $S_8$ )		85.35	
$T_{26} \rightarrow A_{20}$		91.44	
$T_{26} \rightarrow A_{20} \rightarrow A_{14}$	90.45 $\pm$ 0.122	90.66 $\pm$ 0.120	<b>90.81 <math>\pm</math> 0.124</b>
$T_{26} \rightarrow A_{20} \rightarrow A_{14} \rightarrow S_8$	86.71 $\pm$ 0.163	86.85 $\pm$ 0.250	<b>86.97 <math>\pm</math> 0.230</b>

### B. Teacher Selection Ablation

In this section, we investigate the effect of different method of teacher selection. For this we distilled a compact student network (4-layered CNN) from a teacher (10-layered CNN) on CINIC-10 dataset. We performed three types of teacher selection: (1) baseline selection in which we perform our proposed CES-KD method. Mainly, the bucketed sorted data are assigned as following: the easiest samples to the lowest capacity teacher assistant and the hardest examples to the largest teacher network; (2) Anti-selection in which we assign the easy buckets to the largest teacher network and the buckets containing the hardest examples are assigned to the lowest capacity teacher assistant network. Finally, (3) random curriculum in which we randomly assign teacher assistants and teachers to the buckets during training.

TABLE IV

TEST ACCURACY USING PLAIN CNN ARCHITECTURE ON CINIC-10. WE CONSIDERED THE DISTILLATION PATH  $T_{10} \rightarrow T_8 \rightarrow T_6 \rightarrow T_4$ . WE PERFORMED DIFFERENT SCENARIOS OF TEACHER SELECTION: BASELINE, ANTI, AND RANDOM SELECTION.

Selection	Baseline	Anti	Random
Accuracy	<b>71.635 <math>\pm</math> 0.179</b>	71.002 $\pm$ 0.172	71.156 $\pm$ 0.036

Table IV shows the test accuracy on CINIC-10 of a student network (4-layered plain CNN). We demonstrate that having a less specialized expert on an easy concept samples and a very specialized expert on hard samples enhances the performance of a student model. The anti-selection gave the least performance overall. This validates our previous hypothesis that larger networks might not be the optimal teacher when it comes to teaching simple concepts through simple examples. As for the random curriculum where at each distillation step a random expert is chosen, this just shows that there exists a selection rule that can enhance the performance of the distillation process of a student network.

### C. CES-KD for faster student training

TAKD and DGKD methods play a distinctive role in bridging the capacity gap problem. However, for edge devices where the training or fine-tuning time is limited, performing ensemble guidance to the training of the student can be impractical to implement. In Figure 4 we show the training loss and the test accuracy curves over epochs of a 4 layer CNN model and ResNet8 student implemented for all TAKD, DGKD, and

TABLE V

TEST ACCURACY (%) OF STUDENT NETWORKS ON CIFAR100 ON DIFFERENT STATE-OF-THE-ART DISTILLATION METHODS USING DIFFERENT NETWORK ARCHITECTURES. (\*) ARE VALUES PROVIDED IN [28] AND (\*\*) ARE IMPLEMENTED VALUES FROM PUBLICLY AVAILABLE CODE. AVERAGE OVER 5 RUNS.

Teacher Student	NOKD (*)	BLKD (*)	FitNet (*)	AT (*)	SP (*)	CC (*)	VID (*)	RKD (*)	PKT (*)	AB (*)	FT (*)	FSP (*)	NST (*)	CRD (*)	TAKD (**)	DGKD (**)	CES-KD (ours)
wrn-40-2	75.61	74.92	73.58	74.08	73.83	73.56	74.11	73.35	74.54	72.50	73.25	72.91	73.68	<u>75.48</u>	75.25	<u>75.38</u>	<b>75.70</b>
wrn-16-2	73.26																
ResNet110	74.31	70.67	68.99	70.22	70.04	69.48	70.16	69.25	70.25	69.53	70.22	70.11	69.53	<u>71.46</u>	71.15	<u>71.48</u>	<b>71.59</b>
ResNet20	69.06																

CES-KD (ours) method. We see that our method converges faster than DGKD and TAKD. Indeed, we see a fast drop in the training loss with CES-KD using both architectures 4-layered CNNs and ResNet8, after epoch 30 when compared to TAKD and DGKD. Our technique leads to higher performance in fewer epochs. In fact, to reach a target test accuracy of 61%, our method needs only 30 epochs whereas TAKD and DGKD reach this target value only by epoch 90 for the plain CNN architecture on CIFAR-100. Similarly, for the ResNet8 architecture, if a target test accuracy is set at around 86%, the CES-KD method reaches this value by epoch 30 while TAKD and DGKD take much more epochs to reach it. This rapid convergence of the training curve of the student is linked to the curriculum on both data and teacher assignment. This was previously observed in other works in CL [31].

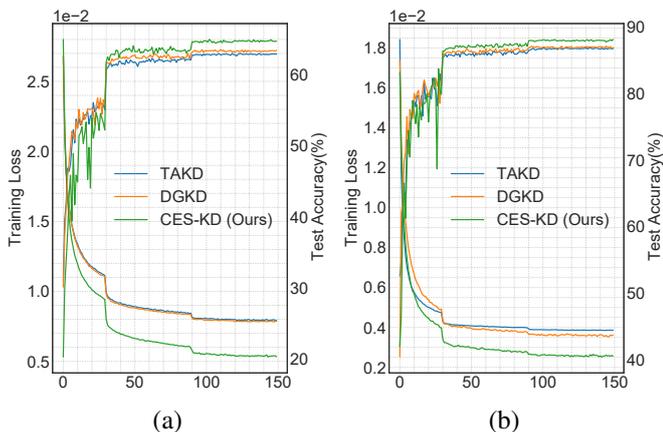


Fig. 4. Training loss and Test Accuracy of a compact student network 4-layered CNN network on CIFAR-100 test set (a) and ResNet8 on CIFAR-100 test set (b) using different TA-based KD methods TAKD, DGKD, and CES-KD (ours).

#### D. CES-KD vs SOTAs in KD

To assess the effectiveness of our method, we compare its test performance to current state-of-the-art KD methods, such as BLKD [7], FitNet [13], AT [32], SP [14], CC [33], VID [34], RKD [35], PKT [36], AB [37], FT [38], FSP [39], NST [40], CRD [41], TAKD [9], and DGKD [10]. In Table V, we provide the test accuracies on CIFAR-100 dataset and using two different architectures WideResNets and ResNets. The bold values are the methods that are outperforming, the single underlined values are the second best and the double underlined values show the third best overall. For TA-based methods (i.e. TAKD, DGKD, and CES-KD), we adopt the following distillation paths for each architecture network: (1) WRN  $40 \times 2(T) \rightarrow$  WRN  $34 \times 2(A_1) \rightarrow$  WRN  $22 \times 2(A_2) \rightarrow$

WRN  $16 \times 2(S)$ , and (2) ResNet110  $\rightarrow$  ResNet56  $\rightarrow$  ResNet44  $\rightarrow$  ResNet32  $\rightarrow$  ResNet20. We observe globally that our method substantially outperforms some of these KD techniques. Mainly we witness an improvement of 2.44% from the student network trained individually and solely from data. Similarly, we see the similar trend for the ResNet architecture having a test accuracy on CIFAR-100 of 71.59%.

#### E. Results on ImageNet

In order to assess the scalability of our method to large datasets, we apply CES-KD to the ImageNet dataset. We chose ResNet56 as our teacher and ResNet18 as our student. The distillation path is ResNet56 (T)  $\rightarrow$  ResNet34 ( $A_1$ )  $\rightarrow$  ResNet18 (S). Table VI shows top1% and top5% validation accuracy on ImageNet. Our method achieved 69.96% top 1% accuracy compared to the baseline vanilla KD, BLKD, which reaches 68.94%. This demonstrates the scalability of our method to larger datasets.

TABLE VI  
TOP 1 ACCURACY (%) ON IMAGENET. THE DISTILLATION PATH IS:  
RESNET34  $\rightarrow$  RESNET26  $\rightarrow$  ResNet18

	Teacher	Student	BLKD (KD)	TAKD	DGKD	CES-KD
Top1%	73.1	68.54	68.94	69.03	69.77	<b>69.96</b>
Top5%	91.2	87.85	87.88	88.21	89.52	<b>89.98</b>

## VI. CONCLUSION

In this paper, we proposed a curriculum guided expert selection for bridging the capacity gap problem in KD. We follow the TA-based KD method but instead of guiding each TA network sequentially or by performing aggregations, we specifically exploit a curriculum on both data and available TA networks to guide the student’s distillation process through a stratified manner of learning. Empirically, we have shown that the learning process of a student through KD is dependent on the difficulty of the samples and also based on the quality of knowledge it is getting from the representative expert. We have demonstrated that a compact student network learns better with lower capacity Teacher/TA networks on easy data samples. Similarly, this compact student’s learning benefits from the knowledge given by larger capacity teacher/TA networks. Our thorough experiments showed that our method can substantially improve the compact student’s performance and that it is comparable to current state-of-the-art and well-performing KD techniques. Finally, we also have shown that our method is scalable to larger datasets. This research was enabled in part by support provided by Compute Canada ([www.computeCanada.ca](http://www.computeCanada.ca)).

## REFERENCES

- [1] S. Han, J. Pool, J. Tran, and W. J. Dally, "Learning both weights and connections for efficient neural networks," 2015.
- [2] Y. Guo, A. Yao, and Y. Chen, "Dynamic network surgery for efficient dnns," *arXiv preprint arXiv:1608.04493*, 2016.
- [3] J. Wu, C. Leng, Y. Wang, Q. Hu, and J. Cheng, "Quantized convolutional neural networks for mobile devices," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4820–4828.
- [4] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, and R. Fergus, "Exploiting linear structure within convolutional networks for efficient evaluation," in *Advances in neural information processing systems*, 2014, pp. 1269–1277.
- [5] P. Luo, Z. Zhu, Z. Liu, X. Wang, and X. Tang, "Face model compression by distilling knowledge from neurons," in *Thirtieth AAAI conference on artificial intelligence*, 2016.
- [6] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," 2018.
- [7] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [8] J. H. Cho and B. Hariharan, "On the efficacy of knowledge distillation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 4794–4802.
- [9] S. I. Mirzadeh, M. Farajtabar, A. Li, N. Levine, A. Matsukawa, and H. Ghasemzadeh, "Improved knowledge distillation via teacher assistant," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 5191–5198.
- [10] W. Son, J. Na, J. Choi, and W. Hwang, "Densely guided knowledge distillation using multiple teacher assistants," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 9395–9404.
- [11] D. Weinshall, G. Cohen, and D. Amir, "Curriculum learning by transfer learning: Theory and experiments with deep networks," in *International Conference on Machine Learning*. PMLR, 2018, pp. 5238–5246.
- [12] C. Bucilua, R. Caruana, and A. Niculescu-Mizil, "Model compression," in *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2006, pp. 535–541.
- [13] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," *arXiv preprint arXiv:1412.6550*, 2014.
- [14] F. Tung and G. Mori, "Similarity-preserving knowledge distillation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1365–1374.
- [15] A. Jafari, M. Rezagholizadeh, P. Sharma, and A. Ghodsi, "Annealing knowledge distillation," *arXiv preprint arXiv:2104.07163*, 2021.
- [16] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proceedings of the 26th annual international conference on machine learning*, 2009, pp. 41–48.
- [17] G. Hacohen and D. Weinshall, "On the power of curriculum learning in training deep networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 2535–2544.
- [18] L. Xiang, G. Ding, and J. Han, "Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification," in *European Conference on Computer Vision*. Springer, 2020, pp. 247–263.
- [19] H. Zhao, X. Sun, J. Dong, Z. Dong, and Q. Li, "Knowledge distillation via instance-level sequence learning," 2021.
- [20] G. Panagiotatos, N. Passalis, A. Iosifidis, M. Gabbouj, and A. Tefas, "Curriculum-based teacher ensemble for robust neural network distillation," in *2019 27th European Signal Processing Conference (EUSIPCO)*, 2019, pp. 1–5.
- [21] X. Wu, E. Dyer, and B. Neyshabur, "When do curricula work?" *CoRR*, vol. abs/2012.03107, 2020. [Online]. Available: <https://arxiv.org/abs/2012.03107>
- [22] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [23] L. N. Darlow, E. J. Crowley, A. Antoniou, and A. J. Storkey, "Cinic-10 is not imagenet or cifar-10," *arXiv preprint arXiv:1810.03505*, 2018.
- [24] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [25] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [26] S. Zagoruyko and N. Komodakis, "Wide residual networks," *arXiv preprint arXiv:1605.07146*, 2016.
- [27] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," 2017.
- [28] Y. Tian, D. Krishnan, and P. Isola, "Contrastive representation distillation," in *International Conference on Learning Representations*, 2020.
- [29] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska *et al.*, "Overcoming catastrophic forgetting in neural networks," *Proceedings of the national academy of sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.
- [30] R. Kemker, M. McClure, A. Abitino, T. Hayes, and C. Kanan, "Measuring catastrophic forgetting in neural networks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [31] X. Wu, E. Dyer, and B. Neyshabur, "When do curricula work?" *arXiv preprint arXiv:2012.03107*, 2020.
- [32] N. Komodakis and S. Zagoruyko, "Paying more attention to attention: improving the performance of convolutional neural networks via attention transfer," in *ICLR*, 2017.
- [33] B. Peng, X. Jin, J. Liu, D. Li, Y. Wu, Y. Liu, S. Zhou, and Z. Zhang, "Correlation congruence for knowledge distillation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5007–5016.
- [34] S. Ahn, S. X. Hu, A. Damianou, N. D. Lawrence, and Z. Dai, "Variational information distillation for knowledge transfer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9163–9171.
- [35] W. Park, D. Kim, Y. Lu, and M. Cho, "Relational knowledge distillation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3967–3976.
- [36] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2536–2544.
- [37] B. Heo, M. Lee, S. Yun, and J. Y. Choi, "Knowledge transfer via distillation of activation boundaries formed by hidden neurons," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 3779–3787.
- [38] J. Kim, S. Park, and N. Kwak, "Paraphrasing complex network: Network compression via factor transfer," *arXiv preprint arXiv:1802.04977*, 2018.
- [39] J. Yim, D. Joo, J. Bae, and J. 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*, 2017, pp. 4133–4141.
- [40] Z. Huang and N. Wang, "Like what you like: Knowledge distill via neuron selectivity transfer," *arXiv preprint arXiv:1707.01219*, 2017.
- [41] Y. Tian, D. Krishnan, and P. Isola, "Contrastive representation distillation," *ArXiv*, vol. abs/1910.10699, 2020.