

Computation-Efficient Knowledge Distillation via Uncertainty-Aware Mixup

Guodong Xu¹, Ziwei Liu², Chen Change Loy²

¹The Chinese University of Hong Kong, ²Nanyang Technological University

xg018@ie.cuhk.edu.hk, {ziwei.liu, ccloy}@ntu.edu.sg

Abstract

Knowledge distillation, which involves extracting the “dark knowledge” from a teacher network to guide the learning of a student network, has emerged as an essential technique for model compression and transfer learning. Unlike previous works that focus on the accuracy of student network, here we study a little-explored but important question, i.e., knowledge distillation efficiency. Our goal is to achieve a performance comparable to conventional knowledge distillation with a lower computation cost during training. We show that the UNCertainty-aware mIXup (UNIX) can serve as a clean yet effective solution. The uncertainty sampling strategy is used to evaluate the informativeness of each training sample. Adaptive mixup is applied to uncertain samples to compact knowledge. We further show that the redundancy of conventional knowledge distillation lies in the excessive learning of easy samples. By combining uncertainty and mixup, our approach reduces the redundancy and makes better use of each query to the teacher network. We validate our approach on CIFAR100 and ImageNet. Notably, with only 79% computation cost, we outperform conventional knowledge distillation on CIFAR100 and achieve a comparable result on ImageNet. The code is available at: <https://github.com/xuguodong03/UNIXKD>.

1. Introduction

“Young children show tremendous versatility in their learning and plasticity ... more efficient, both in how they learn and in how their brain is connected within itself.”

- The Gardener and The Carpenter

In pedagogy, learning efficiency is an important measure to evaluate the learning speed of a learner. A good learner not only learns well at the end but also learns fast in the process. It is preferred if the learner can achieve a comparable or better performance with less effort. Analogous to human learning in real world, machine learning methods [3, 20]

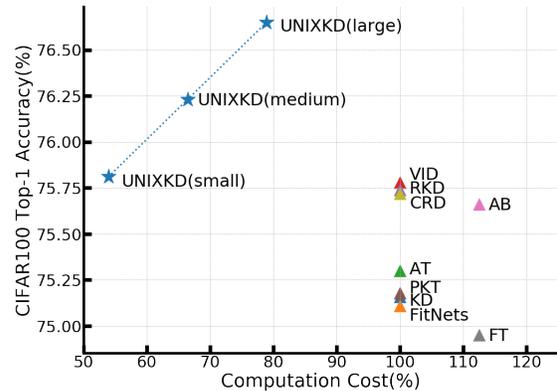


Figure 1. CIFAR100 Top-1 accuracy of different distillation methods. All the methods use the same teacher and student architectures. The computation cost of KD is set as the baseline (100%). Hard label loss is excluded from learning objective. The cost in backbone forward/backward pass is counted and the cost of computing loss is neglected. Our method (UNIXKD) achieves comparable or better performances with a significantly lower computation cost.

also pay attention to learning efficiency issue. The target is to accelerate training from the algorithmic aspect.

Knowledge distillation (KD) [10], serving as a way of model compression [4], has been actively studied to deploy high-performance but cumbersome networks on edge devices. In KD, one typically trains a smaller network (student) under the guidance of a larger network (teacher). The main goal in KD is to obtain a student network with higher accuracy than that trained from scratch. Many studies pursue this goal by transferring better knowledge, such as intermediate features [8, 9, 24, 37, 39] and similarities between samples [22, 32, 33, 35]. The student performance improves rapidly as more powerful algorithms emerge.

Existing works pay relatively more attention to the accuracy dimension but explore little on the aspect of knowledge distillation efficiency. The study of learning efficiency is not only practically beneficial for edge computing and budget-limited training, but also scientifically meaningful. For instance, it reveals the redundancy of the original method and sheds light on factors that influence the performance the most. The study about learning efficiency in KD is

scarce. In this work, we take the first step towards this little-explored but important question.

The challenge in KD efficiency lies in two aspects, *i.e.*, how to define efficiency quantitatively and how to improve efficiency. Previous works in other deep learning topics usually measure efficiency from time dimension, *e.g.*, number of iterations or epochs. However, different distillation methods introduce different operations into the conventional teacher-student framework, hindering a direct comparison on the cost of each iteration across the different methods. Hence, we propose to define efficiency from the standpoint of computation cost, *i.e.*, the number of forward/backward passes in teacher and student networks. With this definition, different KD methods can be compared directly. A more efficient method can achieve the same student accuracy with less computation.

Based on this definition, we propose to improve efficiency by combining two seemingly orthogonal directions, *i.e.*, select informative samples and compact knowledge. The main idea is to measure the uncertainty of each sample in the student forward pass and then compact the knowledge via adaptive mixup according to their level of uncertainty. Specifically, for the uncertain (informative) samples, a mild mixup is applied. For the certain (less informative) samples, a heavy mixup is applied. All the mixed images are passed as queries to teacher. Teacher’s output logits are treated as labels to supervise the learning of the student.

The importance of each training sample is not equal [13]. In the training stage, a student network has different masteries of different samples. The repetition of mastered samples occupies computation but brings nothing new to the network. We adopt uncertainty to estimate the mastery of each sample. The occurrence frequency of samples with high confidence is reduced, so that the network can focus on those uncertain samples. Interestingly, estimating mastery via uncertainty is closely related to active learning (AL). To save annotation cost, AL selects only an informative subset from the original dataset to query the oracle (expert annotators). The KD framework is naturally similar to AL. The teacher plays the role of oracle and is responsible for providing labels. A time-consuming teacher forward pass corresponds to the label expensiveness in AL. Hence, we believe that the classic sampling strategy in AL, *i.e.*, uncertainty-based sampling, is a viable strategy to decide the importance of each sample.

After obtaining the uncertainty of each sample, we compact knowledge through mixup [40] operation. Mixup was originally proposed for data augmentation to improve generalization ability. It compacts the content in two images into a single image via pixel-wise convex combination. Though it can compact knowledge, a disadvantage brought by mixup is that the objects in two images occlude each other, making both of them hard to recognize. Hence, we

adopt adaptive mixup according to uncertainty. A mild mixup is applied to uncertain images so that the important information in them will not be hurt by the mixing image.

Contributions. We take the first step towards a little-explored but important question, *i.e.*, KD efficiency. We quantitatively define the efficiency from the standpoint of forward/backward computation cost, allowing a direct comparison across different KD methods. We propose a novel framework called UNIXKD to improve KD efficiency and conduct thorough experiments to validate the method. With a significantly lower computation cost, we achieve a comparable and even better performance, compared to the conventional KD. We also conduct careful analyses to facilitate a deeper understanding of KD. We show that the redundancy of KD lies in excessive learning of easy categories, and our methods can effectively reduce the redundancy.

2. Related Work

Knowledge Distillation. The goal of KD is to transfer the knowledge from a large teacher network to a small student network. Hinton *et al.* [10] propose to match the category distributions of the teacher and student models via KL-divergence. The relative probability assigned to incorrect categories encodes semantic similarity between similar categories. This “dark knowledge” [10] is shown to benefit the student’s learning. Subsequent studies continue to improve the transfer of knowledge via different learning objectives. Some works propose to mimic the intermediate feature [24] or feature’s transformed variants [9, 12, 14, 36, 37, 39]. Other works [2, 19, 22, 32, 35] extend the notion of unary matching to binary relation mimicking. Dabouei *et al.* [5] study the effects of data augmentation on KD and demonstrates that augmented images can transfer extra knowledge. All the aforementioned methods focus only on improving the accuracy of student.

A recent work by Wang *et al.* [34] highlight the efficiency issue. They also consider both active learning and mixup to tackle the problem. However, our method differs significantly in the following aspects. First, the meaning of efficiency in our work is more holistic. Their goal is only to reduce the query numbers of teacher model, but the student cost is ignored. In fact, if there are 20k initial unlabeled images, their method would require 2 billion student forward passes to select suitable images to query, which is not a negligible expense. In contrast, we consider all the computation cost in the whole KD process. The synergy of mixup and active learning in the two works are also different. They start from a small number of unlabeled images and use mixup to enlarge the dataset to avoid overfitting, while we use mixup to compact multiple images into one. Finally, instead of employing active learning to reduce the billion-level image numbers into a tractable amount, we employ uncertainty to

estimate the value of each original image for the subsequent adaptive mixup.

Several other methods [15, 18, 21] study the data efficiency problem in KD with the objective of obtaining high student accuracy in a data-limited setting. These methods typically involve synthesizing a vast number of images and using them to perform conventional KD. In contrast, we focus on training efficiency by conducting a holistic analysis on computation cost of student learning irrespective of the amount of labeled data.

Active Learning. The idea behind active learning is that a learning agent can achieve greater accuracy with few training labels if it is allowed to select data by itself. An active agent poses queries in the form of unlabeled instances to be labeled by an oracle (e.g., expert annotator). Active learning methods work well in the setting where labels are expensive. Since obtaining label in KD is also an expensive process, we adopt the strategy in active learning to help select informative samples. Common query strategies include uncertainty sampling [17, 26, 29], query-by-committee [28] and expectation-based methods [27, 25]. We adopt uncertainty-based methods as the sampling strategy due to its low computation expense.

3. Methodology

The conventional knowledge distillation [10] minimizes the distance between teacher and student softened logits through KL-divergence. Subsequently, numerous methods explore distilling various other kinds of knowledge, such as feature map [24], attention map [39], pre-activation [9] and batchnorm statistics [36]. All the distillation methods can be written in a general format:

$$L = d(T_t(G_t), T_s(G_s)), \quad (1)$$

where t and s denote the teacher and student, respectively, G_* denotes the feature of some network, T_* denotes some pre-defined transformation, d is a distance metric. It aims at aligning the transformed teacher features $T_t(G_t)$ and transformed student features $T_s(G_s)$ in the distance space d .

3.1. Computation Cost in KD

As shown in Fig. 2, the computation cost in KD arises from three aspects, i.e., teacher forward pass, student forward pass and student backward pass. Let F_t , F_s and B_s denote the floating point operation numbers in these three passes and N be the batch size. The common practice in KD is that all the training samples would go through all the three passes. So the total computation cost in each iteration is:

$$E_{kd} = N \cdot (F_t + F_s + B_s). \quad (2)$$

A more fine-grained cost is:

$$E_{kd} = N_t \cdot F_t + N_{s1} \cdot F_s + N_{s2} \cdot B_s, \quad (3)$$

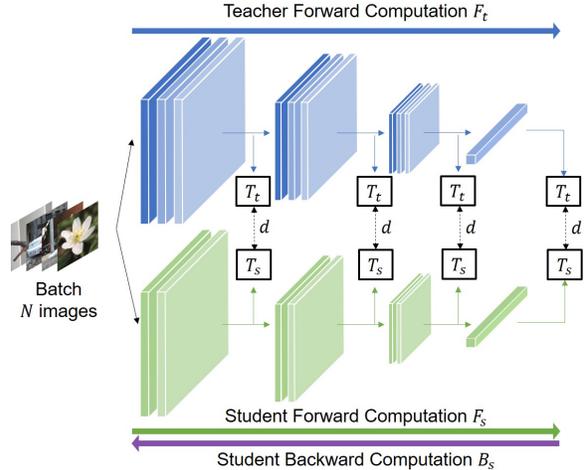


Figure 2. The general framework of distillation. The computation cost consists of three parts, i.e., teacher forward cost F_t , student forward cost F_s and student backward cost B_s .

where N_t , N_{s1} and N_{s2} are the number of samples that go through aforementioned three passes. However, not all samples need to go through all the three passes. An efficient KD method can seek a trade-off between the three terms. Notice that $F_t \gg F_s \approx B_s$. We can try to reduce the teacher query numbers N_t at the expense of an increase in N_{s1} , but the latter is only minor therefore resulting in a net reduction in the total computation cost. On the contrary, previous work [34] reduces N_t into an acceptable range but puts no attention to the hefty increase in N_{s1} , resulting in a net total cost that is intractable.

3.2. Uncertainty-Aware Mixup

The computation of a teacher network is usually much larger than that of the student for the same input image. Hence, to save the total computation, we propose to reduce the number of images fed to the teacher at the expense of a slight computation increase at the student end. Specifically, as shown in Fig. 3, we first feed all the images in a batch to the student network and obtain their uncertainties. We then sort images based on their uncertainty in a descending order and select top k of them to apply adaptive mixup. Finally, the mixed images are fed into two networks to perform the KD process. Since our method only affects the way to feed data and does not change the learning objective, it can be combined with any distillation methods.

Uncertainty Estimation. The importance of each training sample is not equal [13]. Considering the cost of querying teacher, we propose to feed the most informative samples to teacher network. We adopt the uncertainty to measure the informativeness of each training sample. Intuitively, a sample that cannot be classified by the student confidently is a hard sample and it can bring the model with more information upon query. We use entropy to measure the uncertainty

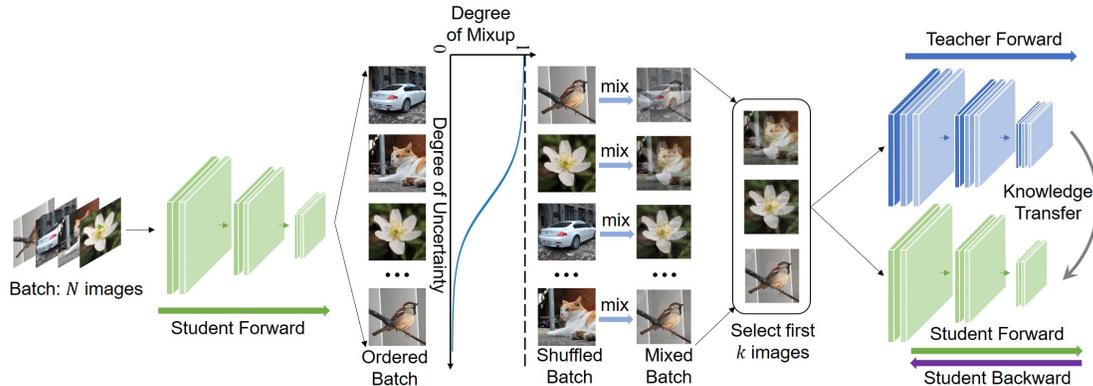


Figure 3. The framework of UNIX for knowledge distillation. Images in a batch are fed into student network and then are sorted in an uncertainty descending order. Mixup is applied adaptively according to uncertainty. For the uncertain samples, a mild mixup is applied and vice versa. Top k mixed images are selected to perform a conventional KD process. Compared to conventional KD, our method saves much cost at the teacher end but increases little cost at the student end. The total cost is reduced due to the asymmetrical capacity between teacher and student.

of classification result:

$$U_c(x) = -\sum_{i=1}^C p_i(x) \log p_i(x), \quad (4)$$

where $p_i(x)$ is the probability that sample x belongs to class i , and C is the number of class. A large entropy indicates a confused classification result. Teacher’s feedback of these samples would benefit the student most. On the contrary, samples that the student has mastered are less useful and can be discarded.

Another two common uncertainty criteria are confidence-based and margin-based measures:

$$U_c(x) = -\max_i p_i(x), \quad (5)$$

$$U_m(x) = -(p_i(x) - p_j(x)), \quad (6)$$

where i and j in Eq. (6) are the first and second most probable class labels, respectively, under the model. Confidence-based criterion considers the largest probability. A large U_c means the model does not assign a high probability to any category, thus the sample is an uncertain sample. Margin-based criterion corrects the shortcoming in confidence-based one, by incorporating the posterior of the second most likely label. Intuitively, instances with large margins are easy, since the classifier has little doubt in differentiating between the two most probable class labels. We will compare the three uncertainty criteria in Sec. 4.2.

Uncertainty-Aware Mixup. Mixup is originally a multi-sample augmentation strategy. It merges each image with another image in a pixel-wise manner:

$$\lambda \sim \text{Beta}(\alpha, \alpha), \quad (7)$$

$$x = \lambda x_i + (1 - \lambda)x_j, \quad (8)$$

where x_i and x_j are the i -th and j -th image, respectively, and a merging coefficient λ is sampled from a Beta distribution parameterized by α . In original Mixup [40], the x ’s label is a convex combination of x_i ’s label and x_j ’s label. However, in this work, we feed the mixed images to teacher network and treat teacher’s logits as the groundtruth label.

The goal of mixup in our work is not to augment the original data, but to compress the content in two images into a single image. We hope that a mixed image can transfer more knowledge than a normal image. However, compression usually brings information loss. A pixel-wise merging leads to mutual aliasing between two images. The visual pattern in each image is hurt by the mixing image, making the synthesized image blurry and semantically meaningless.

To make use of the compression effect while reducing the destruction to informative samples, we apply mixup operation in an adaptive manner. Specifically, we sort the samples in a batch in a descending order according to their uncertainties (Fig. 3). We also shuffle the original batch to prepare the mixing data. Before the mixup between sorted version and shuffled version, we introduce a correction factor c to control the mixup level:

$$x = (1 - c \cdot \lambda)x_{\text{sort}} + c \cdot \lambda x_{\text{shuffle}}, \quad (9)$$

where x_{sort} and x_{shuffle} are both 4D tensors with the first dimension equals to the batchsize, the correction factor c is a real number that belongs to interval $[0, 1]$. The c ’s value is related to the uncertainty ranking of each sample. For an uncertain sample, c is set to a small number and $1 - c \cdot \lambda$ is close to 1, so that mixup is mild and vice versa. In fact, c can be any monotonic increasing function of uncertainty ranking. In this work, we select the sigmoid function:

$$c = \text{sigmoid}\left(w \cdot \frac{\text{ranking} - b}{\text{batchsize}}\right), \quad (10)$$

Table 1. Results on cross-architecture pairs: CIFAR100 Top-1 accuracy and computation cost. **Bold** font denotes the result that outperforms KD. Our method averagely surpasses KD by 1.07% with 79.44% computation cost. Average on 4 runs.

Teacher	Student	VGG13 MobileNetV2	ResNet50 MobileNetV2	ResNet50 VGG8	resnet32×4 ShuffleV1	resnet32×4 ShuffleV2	WRN-40-2 ShuffleV1
F_t/F_s		38.17	174.00	13.56	27.16	23.49	8.22
KD	Acc	68.26	68.26	73.46	74.45	75.16	75.78
	Computation	100%	100%	100%	100%	100%	100%
Random+KD $k=48$	Acc	66.99	65.72	73.32	73.69	74.68	75.70
	Computation	75%	75%	75%	75%	75%	75%
UNIXKD $k=40$	Acc	67.46	67.78	73.62	75.44	76.23	76.59
	Computation	64.99%	63.07%	68.93%	65.93%	66.42%	72.29%
UNIXKD $k=48$	Acc	68.47	69.06	74.24	76.41	76.65	76.92
	Computation	77.49%	75.57%	81.43%	78.43%	78.92%	84.79%

where the b and w are parameters to control the position and shape of sigmoid function. As w increases from small to large, the sigmoid function changes from a linear function to a step function. Finally, we select the top k images from all the mixed images and feed them into the teacher and student networks. The distillation loss in Eq. (1) is computed on these mixed images.

Computation Cost Analysis. The computation in our method comes from two parts, *i.e.*, uncertainty estimation and distillation of mixed images. The total computation cost in each iteration is:

$$\begin{aligned}
 E &= N \cdot F_s + k \cdot (F_t + F_s + B_s) \\
 &= k \cdot F_t + (N + k) \cdot F_s + k \cdot B_s
 \end{aligned}
 \tag{11}$$

Compared to conventional KD in Eq. (2), our method increases the number of student forward passes from N to $N + k$ and reduces the teacher forward passes from N to k . Considering F_t is usually one order larger than F_s , the total cost is reduced.

4. Experiments

The experiments section consists of three parts. In Sec. 4.1, we apply UNIX in conventional KD and demonstrate its ability to save computation. Ablation studies are conducted in Sec. 4.2 to examine the effectiveness of each designed component. In Sec. 4.3, we conduct thorough analyses to understand how our method works.

Evaluations are conducted on CIFAR100 [16] and ImageNet [6] datasets, both of which are widely used as the benchmarks for KD. CIFAR100 consists of 60,000 32×32 colour images, including 50,000 images for training and 10,000 images for testing. There are 100 classes, each contains 600 images. ImageNet is a large-scale classification dataset, containing 1,281,167 images for training and 50,000 images for testing.

Since there is no method designed for improving effi-

ciency in KD¹, we compare our method with two baselines, *i.e.*, conventional KD and Random. Our goal is to achieve a comparable performance with conventional KD with a smaller computation cost. Hence, we regard KD as the upper bound. To demonstrate the improvement brought by our method, we select the Random, *i.e.*, randomly select k images from a batch to perform KD, as the lower bound. We set batchsize as 64 and 256 in CIFAR100 and ImageNet experiments, respectively. We select VGG [30], ResNet [7], WRN [38], ShuffleNet [31, 41] and MobileNetV2 [11] as teacher and student networks.

4.1. Comparative Results

CIFAR100. We compare performances on 11 teacher-student pairs to eliminate architecture influence. We split them into two groups according to whether teacher and student have similar architecture styles. The results are shown in Table 1 and Table 2. In each table, we list F_t/F_s , the computation cost ratio, to facilitate understanding about the capacity gap between teacher and student. For all the methods, we list their accuracy and computation. We regard KD’s computation as the baseline (100%). To study the efficiency in the transferring process, we *exclude the hard label loss* in conventional KD for all the methods. Our method is still advantageous after adding the hard label loss. The results are provided in the Sec. B.

For teacher-student pairs with different architectures (Table 1), our method ($k = 40$, the number of top images from all the mixed images, see Eq. (11)) outperforms KD on four out of six pairs. On average, it surpasses KD by 0.29% with 66.94% computation cost. When we increase k from 40 to 48, our method outperforms KD on all the pairs. On average, it surpasses KD by 1.07% with 79.44% computation cost. With a similar architecture setting (Table 2), our method ($k = 36$) slightly fall behind KD by 0.11% accuracy

¹Although Wang *et al.* [34] involve efficiency, they require the dataset to be a small image pool instead of the whole dataset in conventional KD, making the comparison with our method not suitable.

Table 2. Results on similar-architecture pairs: CIFAR100 Top-1 accuracy and computation cost. **Bold** font denotes the result that outperforms KD. Our method averagely surpasses KD by 0.49% with 92.40% computation cost. The marginal superiority is due to the small gap between teacher and student computation cost. Average on 4 runs.

Teacher	Student	WRN-40-2	WRN-40-2	resnet56	resnet32×4	VGG13	VGG16
		WRN-16-2	WRN-40-1	resnet20	resnet8×4	VGG8	VGG8
F_t/F_s		3.25	3.93	3.06	6.12	2.97	4.14
KD	Acc	75.06	73.95	70.94	73.54	73.50	72.34
	Computation	100%	100%	100%	100%	100%	100%
Random+KD $k=48$	Acc	74.73	73.55	70.09	72.87	72.62	71.68
	Computation	75%	75%	75%	75%	75%	75%
UNIXKD $k=36$	Acc	75.19	73.51	70.06	74.26	73.18	72.38
	Computation	75.31%	73.13%	76.01%	68.57%	76.35%	72.53%
UNIXKD $k=48$	Acc	75.69	74.72	70.77	74.67	73.56	72.77
	Computation	94.06%	91.88%	94.76%	87.32%	95.10%	91.29%

Table 3. CIFAR100 Top-1 accuracy and computation cost when applying UNIX to AT [39] and SP [33]. **Bold** font denotes the result that outperforms original distillation method. On both similar and different architecture settings, UNIX improves the efficiency of AT and SP.

Teacher	Student	F_t/F_s	KD+AT	Random+KD+AT	UNIXKD+AT	KD+SP	Random+KD+SP	UNIXKD+SP
WRN-40-2	WRN-16-2	3.25	75.16 100%	74.85 75%	75.29 75.31%	74.64 100%	73.67 75%	74.66 75.31%
resnet32×4	ShuffleV2	23.49	75.70 100%	74.73 78.13%	76.50 78.92%	75.74 100%	75.05 78.13%	77.12 78.92%

Table 4. ImageNet Top- k accuracy and computation cost. UNIXKD outperforms KD in label-based setting and achieves comparable results in label-free setting. We use a batchsize of 256. We set $k=200$ for Random and $k=192$ for UNIXKD.

	KD	Random+KD	UNIXKD	KD+label	UNIXKD+label
Top-1	45.16	44.28	45.11	53.53	53.64
Top-5	72.57	71.79	72.73	77.75	78.39
Comp.	100%	78.13%	77.23%	100%	77.23%

with 73.65% computation. By increasing k to 48, it outperforms KD by 0.49% with 92.40% computation. The difference in the degree of improvement in the two settings lies in the difference in F_t/F_s . In Table 1, the student networks are mostly lightweight models. The computation is dominated by the number of teacher forward passes. However, in Table 2, the computation of teacher and student is roughly of the same magnitude. Therefore, the decrease of teacher forward pass is compensated by the increase at the student end, causing a marginal difference in the total computation. Hence, our method is more applicable in the setting where the computation gap between teacher and student is large.

Combine with Other KD Methods. Besides conventional KD [10], we also combine UNIX with other distillation methods. Specifically, we select AT [39] and SP [33], for

they represent two mainstream methods in distillation, *i.e.*, mimicking intermediate features and mimicking relations between samples, respectively. Since the learning objective in the two methods cannot update the classifier layer, we combine them with conventional KD, namely the final loss is the combination of KL divergence and respective mimicking loss. We also select two teacher-student pairs to examine our approach under similar and different architecture settings. As shown in Table 3, UNIXKD outperforms AT and SP on both similar-architecture and cross-architecture settings, demonstrating the applicability of UNIX to various distillation methods.

ImageNet. Limited by computation resource, we only conduct one teacher-student pair on ImageNet, *i.e.*, ResNet18 as teacher and ShuffleV2×0.5 as student. As shown in Table 4, with just 77.23% computation cost, UNIXKD outperforms KD on both Top-1 and Top-5 accuracies in label-based setting. It also achieves comparable results in label-free setting. The results on ImageNet demonstrate the scalability of UNIXKD to large-scale dataset.

4.2. Ablation Study

We selectively disable certain parts of the whole framework to examine their effects. The experiments are conducted on CIFAR100. Specifically, there are four variants: 1) Uncertainty without Mixup, namely select the top-

Table 5. Ablation study on CIFAR100 dataset. We compare four variants of UNIXKD with KD and Random. We list their mean accuracies on six teacher-student pairs. **Bold** font denotes result that outperforms KD. The two main components in our methods, *i.e.*, uncertainty and mixup, can both achieve comparable performances with KD. The combination leads to further improvements. Average on 4 runs.

Teacher	Student	WRN40-2	resnet32×4	ResNet50	resnet32×4	resnet32×4	WRN40-2	Mean
		WRN16-2	resnet8×4	VGG8	ShuffleV1	ShuffleV2	ShuffleV1	
KD		75.06	73.54	73.46	74.45	75.16	75.78	74.58
Random $k=48$		74.73	72.87	73.32	73.69	74.68	75.70	74.17
Uncertainty (w/o Mixup)	Confidence	75.34	73.79	73.30	73.88	74.93	75.39	74.44
	Margin	75.50	74.09	73.16	73.97	75.20	75.43	74.56
	Entropy	75.24	74.17	73.27	74.27	75.01	75.86	74.64
Mixup w/o Uncertainty		73.59	72.69	73.11	75.93	75.72	76.36	74.57
Non-adaptive Mixup		74.29	73.34	73.63	76.03	76.03	76.12	74.91
UNIXKD	$w = 1$	75.69	74.46	73.97	75.80	76.23	77.09	75.54
	$w = 10$	75.59	74.67	74.24	76.41	76.65	76.92	75.75
	$w = 1000$	74.93	73.96	73.86	76.32	76.47	76.16	75.28

k uncertain samples and feed them to networks directly, 2) Mixup without Uncertainty, *i.e.*, randomly select k images and apply conventional mixup, 3) Non-adaptive Mixup, namely select top- k uncertain samples and apply conventional mixup, 4) UNIXKD, *i.e.*, the full version of our method. For all the variants, we take $k=48$. We list the mean accuracy averaged on six teacher-student pairs to show an overall performance. The results are shown in Table 5.

Effect of Uncertainty Strategy. We eliminate the influence of mixup and examine the improvement brought by uncertainty sampling. As discussed in Sec. 3.2, we compare three sampling strategies, *i.e.*, Confidence, Margin and Entropy. As shown in Table 5, all three uncertainty strategies outperform Random and achieve comparable results with KD. Notably, without the help of mixup, entropy alone can surpass KD. Among these strategies, confidence behaves worse than the other two methods, because it only takes the maximum probability into account and ignores the full distribution in logits.

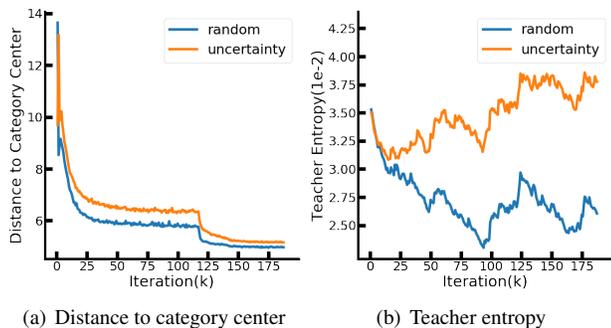


Figure 4. The property of selected samples. Uncertainty strategy prefers to select samples near the decision boundary. These samples are regarded as hard by both teacher and student.

Effect of Mixup. We investigate Mixup’s effect by remov-

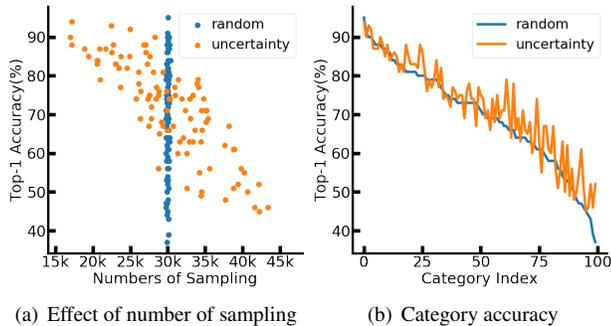
ing the uncertainty estimation. As shown in Table 5, it achieves a similar mean accuracy compared to KD. However, its variance among different pairs is large. It outperforms KD by a large margin on the last three pairs while it falls behind KD on the first three pairs. We conjecture that mixup’s occasional accuracy degradation is a result of the aliasing effect on informative samples. To examine it, we use Uncertainty and Mixup simultaneously while combining them in a non-adaptive manner. The accuracy gap between this variant and the full version shows that adaptive mixup is necessary for mitigating the aliasing effect.

Effect of Correction Factor. As discussed in Sec. 3.2, the degree of adaptive mixup is controlled by the correction factor c , which is the function of uncertainty ranking, and its shape is controlled by the parameter w . We study three cases: 1) $w=1$, c is a linear function of uncertainty ranking, 2) $w=10$, c has a typical sigmoid shape, 3) $w=1000$, c becomes a step function; it is equivalent to applying no mixup to uncertainty samples and applying conventional mixup to certain samples. As shown in Table 5 (bottom), the mean accuracy shows a rise-and-fall trend as w increases from 1 to 1000. And $w=10$ achieves the best results.

4.3. Further Analysis

We further analyze the uncertainty criterion to help understand how it improves accuracy and efficiency. All the experiments are conducted on CIFAR100.

The Property of Selected Samples. We study the sample distribution in feature space by computing the distance between each selected sample and its corresponding category center. We use resnet32×4 and ShuffleV2 as teacher and student networks, respectively. As shown in Fig. 4(a), uncertainty strategy prefers to select samples far away from the category center, *i.e.*, samples near the decision boundary. These samples are more informative than centered sam-

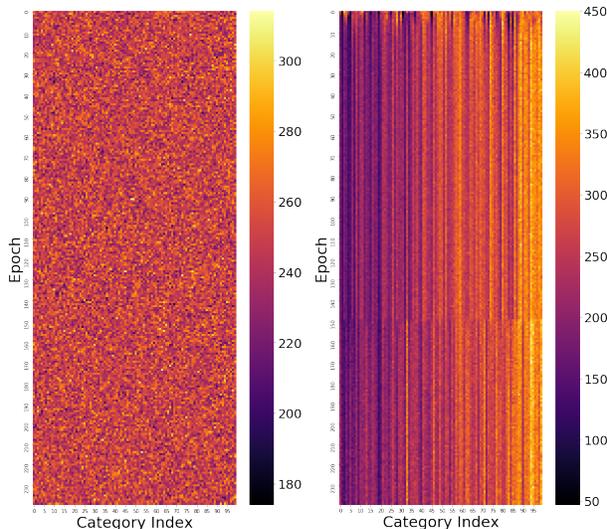


(a) Effect of number of sampling (b) Category accuracy
 Figure 5. Uncertainty strategy selects more samples on hard classes and less samples on easy classes. It greatly improves the performances on hard classes while maintaining comparable performances on easy classes, making the total performance increase.

ples for the student to distinguish different categories. We also compute teacher’s entropy of selected samples as illustrated in Fig. 4(b). For the hard samples selected by student, teacher also regards them as hard, demonstrating that the sample difficulty is an inherent property of data, and an uncertainty strategy can discover informative samples from a nondistinctive dataset. Besides, teacher’s entropy increases as number of iteration increases, showing that Uncertainty’s selection ability evolves as training proceeds.

How Uncertainty Sampling Improves Accuracy? Considering UNIX changes the occurrence frequency of each sample in the whole training stage, we study the relation between the sampling numbers of each category and the corresponding category accuracy.

We use resnet32×4 and resnet8×4 as teacher and student, respectively. The result is shown in Fig. 5(a). Each node in the figure represents a single category. The x -axis denotes the total times that each category has been sampled in the training process. The y -axis denotes the accuracy of each category. For Random baseline, the sampling times of different categories are roughly the same. However, for uncertainty strategy, there is an obvious negative correlation. Uncertainty tends to select more samples on hard categories and fewer samples on easy categories. To have a better view of two methods’ different behaviors on different categories, we arrange the category index in the descending order of accuracy yielded by Random. As shown in Fig. 5(b), Uncertainty achieves higher accuracy on hard categories, benefited from their large sampling numbers. Though uncertainty strategy samples less on easy categories, its accuracy on easy categories remains comparable to that of Random. Hence, the overall accuracy of uncertainty is better than Random. The result suggests that the excessive learning of easy categories causes redundancy in conventional KD. Our method reduces this redundancy by reasonably adjusting the occurrence frequencies of different categories, thus achieving a comparable result with lower computation.



(a) Random sampling pattern (b) Uncertainty sampling pattern
 Figure 6. Random sampling pattern varies little as training proceeds. Uncertainty learns hard category sampling pattern rapidly in the first several epochs, and keeps stable in following epochs.

How Sampling Changes with Epoches? To understand student’s learning behavior along the time dimension, we study how quick the uncertainty criterion is in picking the hard sampling strategy. Specifically, we count the sample numbers of different categories in each epoch and illustrate them through heatmaps in Fig. 6. The category index along x -axis is arranged in a descending order according to the category accuracy of final epoch. Uncertainty’s trend along the category dimension corresponds to hard category sampling. From the time dimension, most of the changes happen at the first several epochs. It indicates that our method learns hard category sampling strategy in a fast speed and then maintains a stable sampling pattern until the end of the training. As a comparison, Random strategy’s pattern varies little along both the time and category dimensions.

5. Conclusion

In this work, we pay attention to a little-explored but important question in knowledge distillation, *i.e.*, KD efficiency. The objective is to reduce the computation cost in KD without compromising its performance. We proposed a novel framework called UNIX to tackle this question. It evaluates the informativeness of each training sample via uncertainty estimation and then applies an adaptive mixup to compact more knowledge into a single sample. We validate our method by conducting thorough experiments on CIFAR100 and ImageNet. Our method achieves better result than conventional KD with less computation, demonstrating the effectiveness of our approach. Further analyses show that our method can effectively reduce the redundancy in knowledge distillation.

References

- [1] <http://tiny-imagenet.herokuapp.com/>. 11
- [2] Sungsoo Ahn, Shell Xu Hu, Andreas Damianou, Neil D. Lawrence, and Zhenwen Dai. Variational information distillation for knowledge transfer. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2
- [3] O. Y. Al-Jarrah, P. D. Yoo, S. Muhaidat, G. K. Karagiannis, and K. Taha. Efficient machine learning for big data: A review. *arXiv preprint arXiv:1503.05296*, 2015. 1
- [4] Cristian Buciluundefined, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '06, page 535–541, New York, NY, USA, 2006. Association for Computing Machinery. 1
- [5] Ali Dabouei, Sobhan Soleymani, Fariborz Taherkhani, and Nasser M. Nasrabadi. Supermix: Supervising the mixing data augmentation. *arXiv preprint arXiv:2003.05034*, 2020. 2
- [6] J. Deng, W. Dong, R. Socher, L. Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 5, 11
- [7] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 5, 11
- [8] Byeongho Heo, Jeesoo Kim, Sangdoon Yun, Hyojin Park, Nojun Kwak, and Jin Young Choi. A comprehensive overhaul of feature distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 1
- [9] Byeongho Heo, Minsik Lee, Sangdoon Yun, and Jin Young Choi. Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In *AAAI*, pages 3779–3787, 2019. 1, 2, 3
- [10] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015. 1, 2, 3, 6
- [11] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. 5, 11
- [12] Seung Hyun Lee, Dae Ha Kim, and Byung Cheol Song. Self-supervised knowledge distillation using singular value decomposition. In *The European Conference on Computer Vision*, 2018. 2
- [13] Angelos Katharopoulos and François Fleuret. Not all samples are created equal: Deep learning with importance sampling. *arXiv preprint arXiv:1803.00942*, 2019. 2, 3
- [14] Jangho Kim, Seonguk Park, and Nojun Kwak. Paraphrasing complex network: Network compression via factor transfer. In *Advances in Neural Information Processing Systems*, pages 2760–2769, 2018. 2
- [15] Akisato Kimura, Zoubin Ghahramani, Koh Takeuchi, Tomoharu Iwata, and Naonori Ueda. Few-shot learning of neural networks from scratch by pseudo example optimization. *arXiv preprint arXiv:1802.03039*, 2018. 3
- [16] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009. 5, 11
- [17] David D. Lewis and William A. Gale. A sequential algorithm for training text classifiers. In Bruce W. Croft and C. J. van Rijsbergen, editors, *SIGIR '94*, pages 3–12, London, 1994. Springer London. 3
- [18] Tianhong Li, Jianguo Li, Zhuang Liu, and Changshui Zhang. Few sample knowledge distillation for efficient network compression. In *cvpr*, June 2020. 3
- [19] Yufan Liu, Jiajiong Cao, Bing Li, Chunfeng Yuan, Weiming Hu, Yangxi Li, and Yunqiang Duan. Knowledge distillation via instance relationship graph. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2
- [20] Roi Livni, S. Shalev-Shwartz, and O. Shamir. On the computational efficiency of training neural networks. In *NIPS*, 2014. 1
- [21] Gaurav Kumar Nayak, Konda Reddy Mopuri, Vaisakh Shaj, R. Venkatesh Babu, and Anirban Chakraborty. Zero-shot knowledge distillation in deep networks. In *icml*, 2019. 3
- [22] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 1, 2
- [23] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017. 11
- [24] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014. 1, 2, 3
- [25] Nicholas Roy and Andrew McCallum. Toward optimal active learning through sampling estimation of error reduction. In *International Conference on Machine Learning*, ICML '01, page 441–448, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc. 3
- [26] Tobias Scheffer, Christian Decomain, and Stefan Wrobel. Active hidden markov models for information extraction. In Frank Hoffmann, David J. Hand, Niall Adams, Douglas Fisher, and Gabriela Guimaraes, editors, *Advances in Intelligent Data Analysis*, pages 309–318, Berlin, Heidelberg, 2001. Springer Berlin Heidelberg. 3
- [27] Burr Settles, Mark Craven, and Soumya Ray. Multiple-instance active learning. In J. C. Platt, D. Koller, Y. Singer, and S. T. Roweis, editors, *Advances in Neural Information Processing Systems*, pages 1289–1296. Curran Associates, Inc., 2008. 3
- [28] H. S. Seung, M. Opper, and H. Sompolinsky. Query by committee. In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, COLT '92, page 287–294, New York, NY, USA, 1992. Association for Computing Machinery. 3
- [29] C. E. Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):379–423, 1948. 3

- [30] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 5, 11
- [31] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V. Le. Mnasnet: Platform-aware neural architecture search for mobile. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 5, 11
- [32] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. In *International Conference on Learning Representations*, 2020. 1, 2, 11
- [33] Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *The IEEE International Conference on Computer Vision*, 2019. 1, 6, 11
- [34] Dongdong Wang, Yandong Li, Liqiang Wang, and Boqing Gong. Neural networks are more productive teachers than human raters: Active mixup for data-efficient knowledge distillation from a blackbox model. In *IEEE Conference on Computer Vision and Pattern Recognition*, June 2020. 2, 3, 5
- [35] Guodong Xu, Ziwei Liu, Xiaoxiao Li, and Chen Change Loy. Knowledge distillation meets self-supervision. In *The European Conference on Computer Vision*, 2020. 1, 2
- [36] Jing Yang, Brais Martinez, Adrian Bulat, and Georgios Tzimiropoulos. Knowledge distillation via adaptive instance normalization. *arXiv preprint arXiv:2003.04289*, 2020. 2, 3
- [37] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 1, 2
- [38] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016. 5, 11
- [39] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *International Conference on Learning Representations*, 2017. 1, 2, 3, 6, 11
- [40] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2018. 2, 4
- [41] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 5, 11

In appendix, we demonstrate the label-free results (training without hard label loss) on TinyImageNet [1] in Sec. A and the label-based (training with hard label loss) results on CIFAR100 [16] in Sec. B. In Sec. C, we provide implementation details of UNIXKD, including network architectures, data augmentation and training hyperparameters.

A. Results on TinyImageNet

We select two teacher-student pairs in TinyImageNet experiments. For all the comparing methods, we exclude the hard label loss in learning objective. As shown in Table 6, UNIXKD outperforms KD by using fewer than 79% of the original computation cost.

B. Label-Based Results on CIFAR100

To study the efficiency in transferring process, we exclude the hard label loss in for the experiments in main paper. We also conduct experiments where hard label loss is incorporated into the learning objective. The results are shown in Table 7 and Table 8. On cross-architecture pairs, our method averagely outperforms KD+label by 1.10% with 79.44% computation cost. On similar-achitecture pairs, our methods averagely outperforms KD+label by 0.53% with 92.40% computation cost. The result shows UNIXKD’s adaptation ability in various settings.

C. Implementation Details

C.1. Network Architectures

On CIFAR100, we adopt resnet [7], ResNet [7], WideResNet [38], MobileNet [11], VGG [30] and ShuffleNet [31, 41] as the network backbones. For resnet, we use resnet-d to represent CIFAR-style resnet with three groups of basic blocks, each with 16, 32 and 64 channels, respectively. resnet8×4 and resnet32×4 indicate a 4× wider network. For ResNet, ResNet-d represents ImageNet-style ResNet with Bottleneck blocks and more channels. For WideResNet (WRN), WRN-d-w represents wide ResNet with depth d and width factor w . For MobileNet, following [32], we use a width multiplier of 0.5. For vgg, ShuffleNetV1 and ShuffleNetV2, we adapt their architectures to CIFAR100 [16] dataset from their original ImageNet [6] counterparts. On ImageNet, we use the PyTorch [23] official implementation of ResNet18 [7] and ShuffleV2×0.5 [31].

C.2. Data Augmentation.

For CIFAR100, we sequentially apply random crop, random horizontal flip and normalization. The crop size is 32×32. For TinyImageNet, we adopt the same operations as CIFAR100, except that the crop size is 64×64. For

Table 6. TinyImageNet Top-1 accuracy and computation cost of different methods. **Bold** font denotes the result that outperforms KD. We use a batchsize of 256 for all the methods.

Teacher Student		resnet32×4 ShuffleV1	resnet32×4 ShuffleV2
F_t/F_s		27.16	23.49
KD	Acc	62.16	62.83
	Computation	100%	100%
Random+KD $k=196$	Acc	62.03	62.07
	Computation	75%	75%
UNIXKD $k=196$	Acc	62.92	63.64
	Computation	78.43%	78.92%

ImageNet, we follow the official implementation of PyTorch [23].

C.3. Training Hyperparameters

CIFAR100. The temperature of KD is set as 4 (same for TinyImageNet and ImageNet). We train all the student models for 240 epochs. The initial learning rate of ShuffleV1, ShuffleV2 and MobileNetV2 is 0.01. For other architectures, the initial learning rate is 0.05. The learning rate is decayed by a factor of 10, respectively, at 150, 180 and 210 epochs. The α of beta distribution is 1.0. We run each experiment on a TITAN-X-Pascal GPU with a batch size of 64. An SGD optimizer with a 5×10^{-4} weight decay and 0.9 momentum is adopted.

When we include hard label loss into the learning objective, we use the common combination coefficient in KD:

$$L = 0.1 \cdot L_{ce} + 0.9 \cdot L_{kd}, \quad (12)$$

where L_{ce} is the cross entropy loss with hard label.

When we combine KD with AT [39] and SP [33], we keep the loss weight of $L_{kd}=1$ and set loss weights of AT and SP to be 1, 000 and 3, 000, respectively.

TinyImageNet. We train all the student models for 100 epochs. The initial learning rate is 0.1 for all the architectures. It is decayed by a factor of 10, respectively, at 40, 70 and 90 epochs. The α of beta distribution is 0.2. We run each experiment on four parallel TITAN-X-Pascal GPU with a total batch size of 256. An SGD optimizer with a 5×10^{-4} weight decay and 0.9 momentum is adopted.

ImageNet. We train all the student models for 90 epochs. The initial learning rate is 0.1 and is decayed by a factor of 10, respectively, at 30 and 60 epochs. The α of beta distribution is 0.2. We run each experiment on eight parallel TITAN-X-Pascal GPU with a total batch size of 256. An SGD optimizer with a 1×10^{-4} weight decay and 0.9 momentum is adopted. When we combine KD with hard label loss, the learning objective is:

$$L = L_{ce} + 0.9 \cdot L_{kd}, \quad (13)$$

Table 7. Results on cross-architecture pairs: CIFAR100 Top-1 accuracy and computation cost. **Bold** font denotes the result that outperforms KD+label. Our method averagely surpasses KD+label by 1.10% with 79.44% computation cost. Average on 4 runs.

Teacher Student	VGG13 MobileNetV2	ResNet50 MobileNetV2	ResNet50 VGG8	resnet32×4 ShuffleV1	resnet32×4 ShuffleV2	WRN-40-2 ShuffleV1	
Teacher Acc	75.38	78.86	78.86	79.58	79.58	76.46	
Student Acc	65.67	65.67	70.68	71.46	72.64	71.46	
F_t/F_s	38.17	174.00	13.56	27.16	23.49	8.22	
KD+label	Acc	68.32	68.24	73.50	74.07	74.93	76.00
	Computation	100%	100%	100%	100%	100%	100%
UNIXKD+label $k=48$	Acc	68.78	68.96	73.88	76.28	76.69	77.04
	Computation	77.49%	75.57%	81.43%	78.43%	78.92%	84.79%

Table 8. Results on similar-architecture pairs: CIFAR100 Top-1 accuracy and computation cost. **Bold** font denotes the result that outperforms KD+label. Our method averagely surpasses KD+label by 0.53% with 92.40% computation cost. Average on 4 runs.

Teacher Student	WRN-40-2 WRN-16-2	WRN-40-2 WRN-40-1	resnet56 resnet20	resnet32×4 resnet8×4	VGG13 VGG8	VGG16 VGG8	
Teacher Acc	76.20	76.20	73.10	79.58	74.73	74.76	
Student Acc	73.45	72.02	69.38	73.22	70.98	70.98	
F_t/F_s	3.25	3.93	3.06	6.12	2.97	4.14	
KD+label	Acc	75.18	74.10	71.06	73.21	73.32	72.22
	Computation	100%	100%	100%	100%	100%	100%
UNIXKD+label $k=48$	Acc	75.48	74.53	70.72	74.78	73.72	73.04
	Computation	94.06%	91.88%	94.76%	87.32%	95.10%	91.29%

where L_{ce} is the cross entropy loss with hard label.