
Improving Ensemble Distillation With Weight Averaging and Diversifying Perturbation

Giung Nam¹ Hyungi Lee¹ Byeongho Heo² Juho Lee^{1,3}

Abstract

Ensembles of deep neural networks have demonstrated superior performance, but their heavy computational cost hinders applying them for resource-limited environments. It motivates distilling knowledge from the ensemble teacher into a smaller student network, and there are two important design choices for this ensemble distillation: 1) how to construct the student network, and 2) what data should be shown during training. In this paper, we propose a weight averaging technique where a student with multiple subnetworks is trained to absorb the functional diversity of ensemble teachers, but then those subnetworks are properly averaged for inference, giving a single student network with no additional inference cost. We also propose a perturbation strategy that seeks inputs from which the diversities of teachers can be better transferred to the student. Combining these two, our method significantly improves upon previous methods on various image classification tasks.

1. Introduction

Deep Ensemble (DE; Lakshminarayanan et al., 2017) averages outputs of multiple models of the same architecture trained with the same data. Despite being simple to implement, DE achieves state-of-the-art performances for various tasks, serving as an oracle for many algorithms to compare against. However, the computational cost of DE scales linearly with the number of models involved in the ensemble, both for training and inference. Especially, the inference cost often becomes critical for a real-world scenario where both memory and time budget is limited.

¹Korea Advanced Institute of Science and Technology, Daejeon, Korea ²Naver, Korea ³AITRICS, Seoul, South Korea. Correspondence to: Giung Nam <giung@kaist.ac.kr>, Juho Lee <juholee@kaist.ac.kr>.

Knowledge Distillation (KD; Hinton et al., 2015) is a method to transfer knowledge from a large teacher network to a smaller student network. The heavy inference cost of DE thus naturally motivates applying KD to reduce it, where a DE is set as a teacher and a single neural network is introduced to be set as a student network. This task of distilling an ensemble teacher network (or distilling from multiple teachers if we treat each ensemble member as a teacher) is often called *ensemble distillation* and has recently been studied actively (Tran et al., 2020; Mariet et al., 2021; Nam et al., 2021; Ryabinin et al., 2021; Du et al., 2020).

For a successful ensemble distillation, one should carefully choose the architecture for a student network. The most straightforward choice would be using a single neural network having the same architecture as the teacher, but this usually yields suboptimal results due to the limited flexibility of the student network. Another choice is to use a student network having *subnetworks*, where the subnetworks share most of the parameters but have a small number of individual parameters, e.g., rank-one factors (Wen et al., 2020; Mariet et al., 2021) or multiple classification heads (Tran et al., 2020). The ensemble distillation with subnetworks are reported to improve performance upon vanilla ensemble distillation (Mariet et al., 2021; Tran et al., 2020; Nam et al., 2021), but they usually require additional computational costs for inference. For instance, the inference cost of a student network having subnetworks defined with rank-one factors scales linearly to the number of subnetworks.

Another important choice for an ensemble distillation algorithm is the training data perturbation strategy. Recently, Nam et al. (2021) studied the importance of diversities in ensemble distillation. When ensemble teachers achieve near-zero train error, the outputs of ensemble teachers for a training data point would be nearly identical, so the ensemble distillation with normal training data would not effectively transfer diversities of the ensemble teachers to students. To this end, Nam et al. (2021) proposed to perturb training data to output-diversifying directions (Tashiro et al., 2020) and use them for distillation. While this indeed improves performance, the perturbation strategy proposed in Nam et al. (2021) only considers teachers without considering how a student would react to such perturbed data. Student

networks are typically less flexible than teacher networks, so a perturbation increasing diversities of teachers may act in an unexpected way when applied to students. Hence, we were motivated to consider both student diversities and teacher diversities into account when designing a perturbation strategy.

In this paper, we propose a novel ensemble distillation method that resolves both of the above-mentioned limitations. Our first contribution is the new way of constructing a student network; we propose to distill with a student with multiple subnetworks during training, but average the subnetwork weights later for inference to get a single student network as a result. As a prototype, we apply this idea to BatchEnsemble (BE; Wen et al., 2020), where the subnetworks are differentiated with multiplicative rank-one factors. Specifically, for such a student network, we propose a training scheme that encourages the subnetwork parameters to stay within the same low-loss region so that the averaged prediction does not degrade the performance while maximally absorbing the diversities transferred from teachers. We show that under a representative training scheme on image classification, averaging those rank-one factors does not degrade predictive performance. The second contribution is a novel perturbation strategy improving upon the one proposed in Nam et al. (2021). Unlike the previous method, our perturbation method considers both students and teachers. More specifically, we find the ‘‘weak points’’ of the student networks by seeking inputs on which the student subnetworks agree with each other (low diversity), and at the same time, find the inputs on which the teachers disagree (high diversity). With this perturbation considering both students and teachers, our method effectively transfers diversities of teachers to students. To demonstrate the effectiveness of our method, we compare ours with previous methods on various image classification benchmarks. We find that ours achieve significantly improved predictive accuracy and uncertainty calibration results without increasing inference cost.

2. Backgrounds

2.1. Setup

The problem we address in this paper is the K -way classification problem; a neural network $\mathcal{F} : \mathbb{R}^D \rightarrow \mathbb{R}^K$ takes D -dimensional inputs \mathbf{x} (i.e., images) and makes predictions about corresponding outputs y (i.e., class label) with K -dimensional logits $\mathcal{F}(\mathbf{x})$. We denote the output probabilities of the model \mathcal{F} for a given input \mathbf{x} as

$$\mathbf{p}_{\mathcal{F}}^{(k)}(\mathbf{x}; \tau) = \frac{\exp(\mathcal{F}^{(k)}(\mathbf{x})/\tau)}{\sum_{j=1}^K \exp(\mathcal{F}^{(j)}(\mathbf{x})/\tau)}, \quad (1)$$

for $k = 1, \dots, K$. Here, we introduce a single scale parameter $\tau > 0$ for temperature scaling which will be used both for training and evaluation procedures.

2.2. Ensemble Distillation

Let $\{\mathcal{T}_1, \dots, \mathcal{T}_M\}$ be a set of pre-trained teachers, and \mathcal{S}_{θ} be a student. In a vanilla ensemble distillation, using the knowledge distillation (KD; Hinton et al., 2015) framework, the student tries to mimic the probabilistic outputs of an ensemble of teachers under the given temperature τ by minimizing the averaged KD loss, which is equivalent to minimizing

$$\tau^2 \mathcal{H} \left[\frac{1}{M} \sum_{m=1}^M \mathbf{p}_{\mathcal{T}_m}(\mathbf{x}; \tau), \mathbf{p}_{\mathcal{S}_{\theta}}(\mathbf{x}; \tau) \right], \quad (2)$$

where $\mathcal{H}[\cdot, \cdot]$ computes the cross-entropy between two probability vectors. Note that this vanilla approach minimizes the discrepancy between the student predictions and the *mean* predictions of the ensemble teachers. As a result, the diversities among ensemble teachers are removed by mean operation and hardly transferred to the student.

2.3. BatchEnsemble and one-to-one distillation

BatchEnsemble (BE; Wen et al., 2020) is a parameter-efficient way to ensemble deep neural networks; each member of the ensemble is constructed in the low-rank subspace with rank-one factors, instead of the full parameter space. With a slight abuse of notation, while a DE would have full set of parameters $\{\theta_1, \dots, \theta_M\}$, BE introduces a shared parameter θ and a set of rank-one matrices $\{\mathbf{r}_1 \mathbf{s}_1^{\top}, \dots, \mathbf{r}_m \mathbf{s}_m^{\top}\}$, and construct m^{th} subnetwork parameter as $\theta \circ \mathbf{r}_m \mathbf{s}_m^{\top}$, where \circ denotes the Hadamard product. Based on BE, Mariet et al. (2021) proposed a one-to-one ensemble distillation scheme, where each BE subnetwork is trying to mimic single ensemble member in one-to-one fashion. Instead of learning the mean prediction of the teachers, the one-to-one distillation minimizes

$$\sum_{m=1}^M \tau^2 \mathcal{H} \left[\mathbf{p}_{\mathcal{T}_m}(\mathbf{x}; \tau), \mathbf{p}_{\mathcal{S}_{\theta \circ \mathbf{r}_m \mathbf{s}_m^{\top}}}(\mathbf{x}; \tau) \right]. \quad (3)$$

That is, each subnetwork copies a member from the ensembles in a one-to-one way. The training procedure for BE ensemble distillation is summarized in Algorithm 1.

2.4. Ensemble distillation with diversifying perturbation

When the ensemble teachers are flexible enough to achieve zero-train error, their responses to a training input would be nearly identical, so a student network distilled from them would not be exposed to the diversities of the teachers. To resolve this, Nam et al. (2021) proposes to perturb training inputs with the Output Diversified Sampling (ODS; Tashiro et al., 2020) that encourages ensemble teachers disagree with each other. The ODS for an input \mathbf{x} is computed as

$$\varepsilon_{\text{ODS}} \propto \nabla_{\mathbf{x}} (\mathbf{w}^{\top} \mathbf{p}_{\mathcal{T}_m}(\mathbf{x}; \tau)), \quad (4)$$

Algorithm 1 Ensemble distillation with BE

Require: Temperature τ , learning rate η .

- 1: **while** not converged **do**
 - 2: Sample an input \mathbf{x} from the train split.
 - 3: **for** $m = 1, \dots, M$ **do**
 - 4: Compute loss for the m^{th} subnetwork:

$$\ell_m \leftarrow \tau^2 \mathcal{H} \left[\mathbf{p}_{\mathcal{T}_m}(\mathbf{x}; \tau), \mathbf{p}_{\mathcal{S}_{\theta \circ (\mathbf{r}_m \mathbf{s}_m^\top)}}(\mathbf{x}; \tau) \right].$$
 - 5: Update rank-one factors:

$$\mathbf{r}_m \leftarrow \mathbf{r}_m - \eta \nabla_{\mathbf{r}_m} \ell_m.$$

$$\mathbf{s}_m \leftarrow \mathbf{s}_m - \eta \nabla_{\mathbf{s}_m} \ell_m.$$
 - 6: **end for**
 - 7: Update shared parameters:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \frac{1}{M} \sum_{m=1}^M \nabla_{\boldsymbol{\theta}} \ell_m.$$
 - 8: **end while**
-

where \mathbf{w} denotes a random guidance vector sampled from the K -dimensional uniform distribution with zero means. Intuitively, the ODS perturbation seeks the direction in the input space to make the output follow the random guidance vector \mathbf{w} , and thus drives ensemble members to produce diverse outputs¹. Nam et al. (2021) demonstrated that a BE student trained with ODS perturbation shows significantly improved performance. They also proposed an improved version of ODS perturbation called ConfODS, where the ODS perturbations are scaled by the confidence of teacher predictions.

3. Improved Ensemble Distillation

3.1. LatentBE : a weight averaged BE student

Utilizing the subnetwork structure, a BE student can capture the diversities of ensemble teachers, but this comes at a cost of increased inference time. To get an output from a BE network, one should execute forward passes M times from scratch since the different rank-one subnetworks do not share hidden layer responses. Hence, although BE significantly reduces the number of parameters compared to DE, its computation cost for inference is identical to that of DE. This is definitely undesirable, especially for distillation where we want a cheap student network applicable for real-world applications. In this section, we propose a novel ensemble distillation framework to circumvent this limitation, where a student network has the same inference cost as a single neural network yet maintains the flexibility to well absorb the diversity of teachers.

Based on the ensemble distillation with BE, we propose a novel framework entitled LatentBE, in a sense that the rank-one factors defining BE are averaged out for inference just as the latent variables are marginalized out for probabilistic

¹Actually, we need a transferability assumption to fully justify this argument. For more detail, please refer to Nam et al. (2021).

Algorithm 2 Ensemble distillation with LatentBE + diversifying perturbation

Require: Temperature τ , learning rate η , rank-one weight

 decay parameter λ , perturbation step size γ .

- 1: **Initialize** rank-one factors to ones:

$$\mathbf{r}_m \leftarrow \mathbf{1} \text{ and } \mathbf{s}_m \leftarrow \mathbf{1} \text{ for } m = 1, \dots, M.$$
 - 2: **while** not converged **do**
 - 3: Sample an input \mathbf{x} from the train split.
 - 4: Sample indices i, j uniformly from $\{1, \dots, M\}$.
 - 5: Perturb the input w.r.t. teacher and student:

$$\tilde{\mathbf{x}} \leftarrow \mathbf{x} + \gamma (\widehat{\text{TDiv}}(\mathbf{x}) - \widehat{\text{SDiv}}(\mathbf{x})).$$
 - 6: **for** $m = 1, \dots, M$ **do**
 - 7: Compute loss for the m^{th} subnetwork:

$$\ell_m \leftarrow \tau^2 \mathcal{H} \left[\mathbf{p}_{\mathcal{T}_m}(\tilde{\mathbf{x}}; \tau), \mathbf{p}_{\mathcal{S}_{\theta \circ (\mathbf{r}_m \mathbf{s}_m^\top)}}(\tilde{\mathbf{x}}; \tau) \right]$$
 - 8: Update rank-one factors:

$$\mathbf{r}_m \leftarrow \mathbf{r}_m - \eta \nabla_{\mathbf{r}_m} \ell_m - \eta \lambda (\mathbf{r}_m - \mathbf{1}).$$

$$\mathbf{s}_m \leftarrow \mathbf{s}_m - \eta \nabla_{\mathbf{s}_m} \ell_m - \eta \lambda (\mathbf{s}_m - \mathbf{1}).$$
 - 9: **end for**
 - 10: Update shared parameters:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \frac{1}{M} \sum_{m=1}^M \nabla_{\boldsymbol{\theta}} \ell_m.$$
 - 11: **end while**
 - 12: **Return the averaged parameter:**

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} \circ \frac{1}{M} \sum_{m=1}^M \mathbf{r}_m \mathbf{s}_m^\top.$$
-

inference. Specifically, we employ a BE student having multiple subnetworks, and do one-to-one ensemble distillation similarly to Algorithm 1, but after training, compute the *weight average* of the rank-one factors to construct a single student network with parameter

$$\boldsymbol{\theta} \circ \left(\frac{1}{M} \sum_{m=1}^M \mathbf{r}_m \mathbf{s}_m^\top \right). \quad (5)$$

After this weight averaging, the inference cost remains the same as that of a single neural network. The idea of weight averaging was first considered in Stochastic Weight Averaging (SWA; Izmailov et al., 2018), where the parameters collected from a single learning trajectory is averaged to construct a better generalizing model. The key observation in SWA is that due to the choice of the specific learning rate scheduling, the parameters in a learning trajectory remain in a wide low-loss region in the loss surface, so averaging them leads to a single robust model. Our LatentBE shares some spirits with SWA but has crucial differences: 1) LatentBE utilizes the diverse subnetworks from the guidance of ensemble teachers, which brings diversity and performance gain of ensemble distillation to *weight average*, and 2) instead of special learning rates schedule of SWA, LatentBE enables *weight average* with the rank-one factors of BE.

The key for making LatentBE successful is, as in SWA, to keep the subnetwork parameters stay in the same low-loss region. For this, 1) we initialized all the rank-one

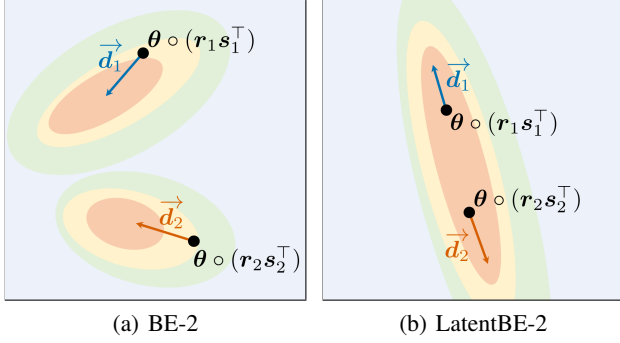


Figure 1. A schematic diagram depicting BE and LatentBE. Here, \vec{d}_1 and \vec{d}_2 denote learning directions obtained from two different teachers. For the BE, (a) two subnetwork parameters *explore* different modes with low-rank subspaces, while for the LatentBE, (b) two subnetwork parameters *expand* a low-rank subspace.

factors $\{r_m, s_m\}_{m=1}^M$ to be one vectors, and 2) set Gaussian prior with one vector mean to those rank-one factors. By doing these, all the rank-one factors start from a similar location and gradually split into individual factors, but they do not deviate too much from each other due to the weight-decay effect driven by the prior. Also, other than these two, we do not require careful learning rate scheduling as in SWA, and this is presumably due to the fact that we are only differentiating rank-one factors for subnetworks with a large body of shared parameters while SWA averages entire parameters.

The rank-one factors for BE and LatentBE play different roles in ensemble distillation. BE aims to find a subspace for each rank-one factor and spread those subspaces to different modes. On the other hand, LatentBE seeks for a flat minima in which all the subspaces defined by the rank-one factors are embedded, and then the rank-one factors are trained to “stretch” the subspace area by letting the subnetworks follow different teacher directions (Figure 1).

3.2. A better perturbation strategy

As we discussed earlier, Nam et al. (2021) utilizes perturbation strategies, ODS and ConfODS, to improve diversity transfer in ensemble distillation. It is an innovative approach for ensemble distillation but does not consider the student networks, especially their diversities. Since there have been several works controlling teacher networks based on the status of student networks for KD (Jin et al., 2019; Mirzadeh et al., 2020), we further conjecture that the ensemble distillation can also be improved with a perturbation strategy considering both teachers and students. Thus, we propose a novel perturbation scheme that considers *both* student and teacher diversities for more effective diversity transfer.

For a given input \mathbf{x} , we measure the functional diversity

of the ensemble of $\{\mathcal{F}_1, \dots, \mathcal{F}_M\}$ by averaging pairwise KL-divergence between probabilistic outputs from different members,

$$\text{Div}(\{\mathcal{F}_m\}_{m=1}^M, \mathbf{x}) = \frac{\sum_{i=1}^M \sum_{j=1}^M D_{ij}(\mathbf{x})}{M(M-1)}, \quad (6)$$

where $D_{ij}(\mathbf{x}) = D_{\text{KL}}(\mathbf{p}_{\mathcal{F}_i}(\mathbf{x}) \parallel \mathbf{p}_{\mathcal{F}_j}(\mathbf{x}))$. From this, we denote the student and teacher diversities as

$$\text{SDiv}(\mathbf{x}) := \text{Div}(\{\mathcal{S}_m\}_{m=1}^M, \mathbf{x}), \quad (7)$$

$$\text{TDiv}(\mathbf{x}) := \text{Div}(\{\mathcal{T}_m\}_{m=1}^M, \mathbf{x}), \quad (8)$$

where $\mathcal{S}_m := \mathcal{S}_{\theta \circ r_m s_m^T}$. We suggest perturbing input to the direction that *minimizes* the student diversity while the teacher diversity is *maximized*,

$$\varepsilon \propto \nabla_{\mathbf{x}} (\text{TDiv}(\mathbf{x}) - \text{SDiv}(\mathbf{x})). \quad (9)$$

The intuition behind this perturbation is as follows: ideally, we want the student subnetworks to learn the diverse outputs of ensemble teachers almost everywhere in the input space. Hence, we first introduce negative student diversity term “ $-\text{SDiv}(\mathbf{x})$ ” to find a point that student subnetworks have low diversities. At the same time, as we originally intended, we use teacher diversity term “ $\text{TDiv}(\mathbf{x})$ ” to improve diversity transfer of ensemble teachers. The combined perturbation thus finds an input point that maximizes the diversity gap between teachers and student subnetworks and gives a strong learning signal to correct it.

In practice, exactly computing (7) and (8) would be costly, so we use stochastic approximations of them where the student and teacher diversities are computed for a randomly selected pair. That is, we pick $i, j \sim \{1, \dots, M\}$ uniformly and compute

$$\widehat{\text{TDiv}}(\mathbf{x}) := D_{\text{KL}}(\mathbf{p}_{\mathcal{T}_i}(\mathbf{x}) \parallel \mathbf{p}_{\mathcal{T}_j}(\mathbf{x})), \quad (10)$$

$$\widehat{\text{SDiv}}(\mathbf{x}) := D_{\text{KL}}(\mathbf{p}_{\mathcal{S}_i}(\mathbf{x}) \parallel \mathbf{p}_{\mathcal{S}_j}(\mathbf{x})), \quad (11)$$

and define the perturbation as

$$\hat{\varepsilon} \propto \nabla_{\mathbf{x}} (\widehat{\text{TDiv}}(\mathbf{x}) - \widehat{\text{SDiv}}(\mathbf{x})). \quad (12)$$

We also found that blocking the gradient flow through one of the teachers or students stabilizes the training,

$$\widehat{\text{TDiv}}(\mathbf{x}) := D_{\text{KL}}(\text{sg}(\mathbf{p}_{\mathcal{T}_i}(\mathbf{x})) \parallel \mathbf{p}_{\mathcal{T}_j}(\mathbf{x})), \quad (13)$$

$$\widehat{\text{SDiv}}(\mathbf{x}) := D_{\text{KL}}(\text{sg}(\mathbf{p}_{\mathcal{S}_i}(\mathbf{x})) \parallel \mathbf{p}_{\mathcal{S}_j}(\mathbf{x})), \quad (14)$$

where $\text{sg}(\cdot)$ denotes the `stop_grad` operation, for example, `.detach()` in PyTorch library.

One thing to note here is that we are directly measuring divergences between teachers or student subnetworks to get

perturbations, unlike the ODS-based perturbation proposed in Nam et al. (2021). An ODS computed from a specific network, in principle, does not guarantee the output diversification of other networks. Hence, the ODS perturbation computed from a single teacher, as suggested in Nam et al. (2021), does not guarantee the diversities among teacher outputs. Nam et al. (2021) argued that this issue can be circumvented by assuming *transferability* of teacher networks, where we assume that the gradients of ensemble teachers are similar to each other, so an ODS computed from a specific teacher generalizes to the other teachers, driving overall diversities among teacher outputs as a result. However, the transferability does not always hold for which an ODS perturbation fails to properly bring diversities. On the other hand, our perturbation directly minimizing or maximizing KL divergences does not require transferability of gradients.

3.3. Improved ensemble distillation algorithm

Our final ensemble distillation algorithm combines two ingredients discussed so far; LatentBE and novel diversifying perturbation. The only overhead during training is the procedure of computing the diversifying perturbation which requires additional forward and backward passes through teacher networks. Our algorithm is summarized in Algorithm 2, with highlights on the part different from the vanilla one-to-one ensemble distillation with BE.

4. Related Works

Ensembles Recent works have shown that an ensemble of deep neural networks can achieve superior performance both in terms of prediction accuracy and uncertainty estimation (Lakshminarayanan et al., 2017; Ovadia et al., 2019). The power of the ensemble comes from *the diversity* among ensemble members, and there have been several works to enhance it, e.g., constructing ensembles with varying hyperparameters (Wenzel et al., 2020), or architectures (Zaidi et al., 2021), or reducing conditional redundancy (Rame & Cord, 2021), or introducing kernelized repulsion (D’Angelo & Fortuin, 2021).

Ensemble distillation The seminal work of Hinton et al. (2015) has already shown the effectiveness of the ensemble distillation. It can be further enhanced by considering the diversity inside the ensemble teacher, e.g., dynamically assign weights to teachers (Du et al., 2020), or treating predictions from teachers as a set of samples from an implicit distribution (Malinin et al., 2020; Ryabinin et al., 2021), or amplifying the diversity via input perturbations (Nam et al., 2021). Besides, several existing approaches propose to use a student having subnetworks which can represent the diversity in predictions (Tran et al., 2020; Mariet et al., 2021; Nam et al., 2021). However, their resulting students

have additional costs for inference and defeat their ends that reduce the computational cost of the ensemble.

5. Experiments

In this section, we present the experimental results on image classification benchmarks including CIFAR-10, CIFAR-100 (Krizhevsky, 2009), TinyImageNet, and ImageNet-1k (Russakovsky et al., 2015). Through the experiments, we empirically validate the following questions:

- How does the subspace discovered by LatentBE look like? - Section 5.1.
- How does the proposed perturbation strategy affect the training of LatentBE? - Section 5.2.
- Does our ensemble distillation algorithm improves performance both in terms of predictive accuracy and uncertainty calibration? - Section 5.3.

Please refer to Appendix B for the training details including data augmentation, learning rate schedules, and other hyperparameter settings.

5.1. Subspaces of LatentBE

In order to investigate the subspaces defined by the subnetwork parameters $\{\theta \circ (r_m s_m^T)\}_{m=1}^M$, we first consider the case of $M = 2$ models where the subspace forms a simple line. More precisely, we parameterize the line passing through two subnetwork parameters as $\{\theta_t \mid t \in \mathbb{R}\}$, where

$$\theta_t = (1 - t) (\theta \circ (r_1 s_1^T)) + t (\theta \circ (r_2 s_2^T)). \quad (15)$$

Figure 2 shows how prediction error and negative log-likelihood vary along the line subspace. The main difference between BE and LatentBE students is the presence of a *loss barrier* between two end-points. The LatentBE student does not have a barrier while the BE student does, as we have depicted in Figure 1. This difference enables the weight averaging of subnetwork parameters for LatentBE. Notably, as shown in (Figure 2, right), the averaged parameter effectively improves negative log-likelihood on the test data, which is consistent with the findings Izmailov et al. (2018) and recent study on the neural network subspaces (Wortsman et al., 2021). These improvements in performance, as we will show later in Section 5.3, confirm the validity of our weight averaging strategy for ensemble distillation.

5.2. Diversification effects of perturbations

Using on the LatentBE-2 model discussed on the previous section, Figure 3 shows the functional diversity (defined in Equation (6)) between two end-points, when distilled with and without ConfODS (Nam et al., 2021) and ours (i.e.,

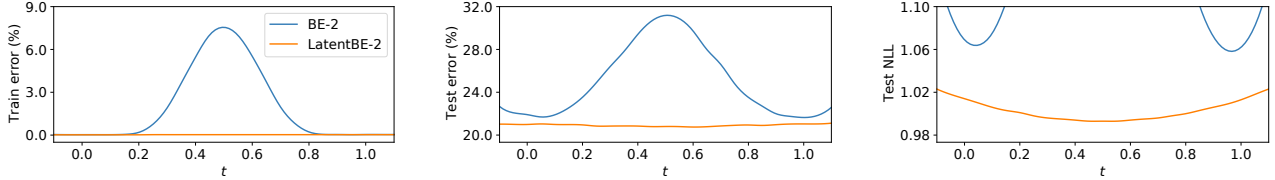


Figure 2. Train errors (left), test errors (middle), and test negative log-likelihood along the line subspace (right) passing through two different subnetwork parameters. t denotes the position on the line as defined in Equation (15). BE-2 and LatentBE-2 students are distilled from the DE-2 teacher for WRN28x4 on CIFAR-100 *without* any input perturbations.

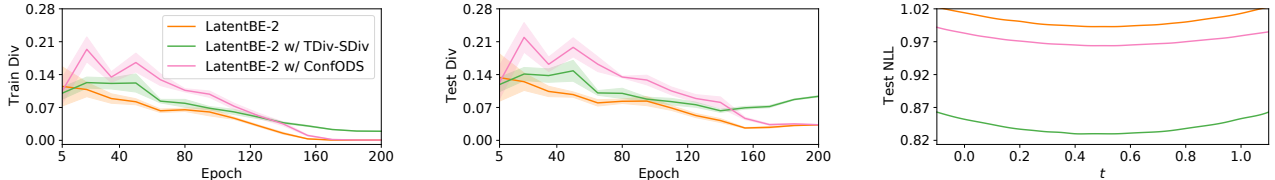


Figure 3. Function diversity in predictions from two subnetwork parameters of the LatentBE-2 student on the train (left) and test (middle) data measured for a training run. We also visualize the resulting student’s test NLL (right). LatentBE-2 students are distilled from the DE-2 teacher for WRN28x4 on CIFAR-100 *with* the stated input perturbations.

TDiv-SDiv). As one can see from the figures, ours better diversifies the subnetworks while maintaining lower test NLL values.

Figure 4 shows the effects of the perturbation strategies during training. For this, we measure the changes in the function diversities of the teacher networks and student subnetworks due to a perturbation ε during a training procedure:

$$\text{TDiv}(\mathbf{x} + \varepsilon) - \text{TDiv}(\mathbf{x}), \quad (16)$$

$$\text{SDiv}(\mathbf{x} + \varepsilon) - \text{SDiv}(\mathbf{x}). \quad (17)$$

Again, our intuition behind the perturbation strategy proposed in Section 3.2 is, pinpointing inputs that should be diversified for students while diversifying learning signals from teachers. In other words, the perturbation that decreases Equation (17) while increasing Equation (16) will be helpful for ensemble distillation.

Figure 4 clearly shows that our proposed perturbation (i.e., TDiv - Sdiv) accomplishes our goal to increase teacher diversities and decrease student diversities. Especially, the student diversities significantly drop during the early stage of training, but gradually increase as the training progress. We conjecture that this is because after enough training the students are diversified for a wide range of inputs, so the effect of $-Sdiv$ perturbation gradually decreases. Although ConfODS exhibits high diversification effects on teachers as intended, it has no significant effect on the students.

As an ablation study, we compared the efficacy of the perturbation strategies in terms of the actual classification performance. Table 1 shows that increasing student diversities (TDiv-SDiv) clearly improves the performance compared to the ones not considering the student diversities (ConfODS or TDiv).

Table 1. Ablation results for perturbation strategies. The results are with WRN28x4 on CIFAR-100.

Method	ACC (\uparrow)	NLL (\downarrow)	ECE (\downarrow)	cNLL (\downarrow)	cECE (\downarrow)
LatentBE-4	79.46 \pm 0.20	0.993 \pm 0.024	0.124 \pm 0.004	0.837 \pm 0.012	0.046 \pm 0.005
+ ConfODS	79.27 \pm 0.30	0.955 \pm 0.008	0.115 \pm 0.002	0.840 \pm 0.007	0.048 \pm 0.004
+ TDiv	79.40 \pm 0.03	0.861 \pm 0.000	0.084 \pm 0.002	0.819 \pm 0.001	0.046 \pm 0.001
+ TDiv-SDiv	80.02\pm0.07	0.792\pm0.004	0.067\pm0.001	0.772\pm0.003	0.041\pm0.003

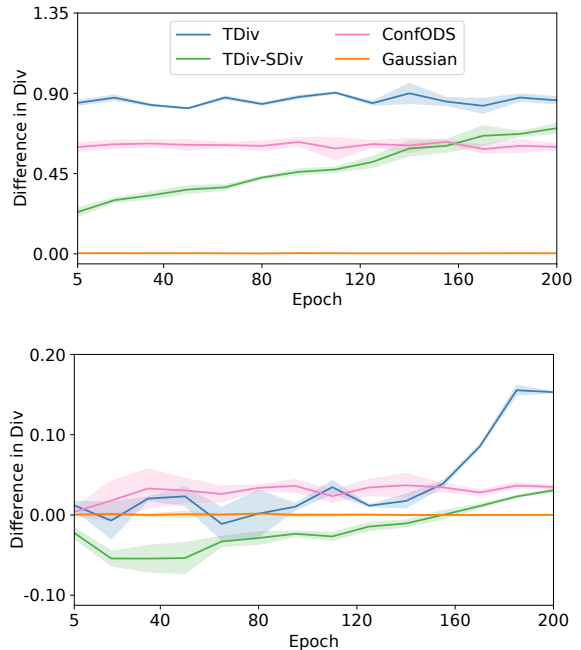


Figure 4. Diversification effects of various perturbations on the train data for the DE teachers (top) and the LatentBE students (bottom). The results are with WRN28x4 on CIFAR-100.

Table 2. Results of the distilled students for WRN28x1 on CIFAR-10. Results with \pm std. are averaged over 4 seeds.

Method	Students distilled from DE-4 teacher					Students distilled from DE-8 teacher				
	ACC (\uparrow)	NLL (\downarrow)	ECE (\downarrow)	cNLL (\downarrow)	cECE (\downarrow)	ACC (\uparrow)	NLL (\downarrow)	ECE (\downarrow)	cNLL (\downarrow)	cECE (\downarrow)
<i>Single forward pass for inference:</i>										
KD (Hinton et al., 2015)	93.70 \pm 0.10	0.270 \pm 0.004	0.042 \pm 0.001	0.201 \pm 0.003	0.011 \pm 0.002	93.68 \pm 0.15	0.273 \pm 0.007	0.042 \pm 0.002	0.202 \pm 0.003	0.010 \pm 0.001
AE-KD (Du et al., 2020)	93.67 \pm 0.05	0.278 \pm 0.002	0.042 \pm 0.000	0.205 \pm 0.001	0.010 \pm 0.001	93.75 \pm 0.16	0.266 \pm 0.008	0.041 \pm 0.001	0.199 \pm 0.004	0.011 \pm 0.001
Proxy-EnD ² (Ryabinin et al., 2021)	93.67 \pm 0.04	0.270 \pm 0.005	0.042 \pm 0.001	0.200 \pm 0.002	0.011 \pm 0.001	94.04 \pm 0.12	0.263 \pm 0.003	0.039 \pm 0.001	0.195 \pm 0.002	0.009 \pm 0.001
KD + LatentBE (Ours)	93.98\pm0.20	0.263\pm0.003	0.041\pm0.002	0.194\pm0.002	0.011\pm0.002	93.97\pm0.10	0.263\pm0.004	0.041\pm0.002	0.194\pm0.002	0.012\pm0.001
+ ConfODS	93.95 \pm 0.12	0.223 \pm 0.007	0.032 \pm 0.001	0.186 \pm 0.004	0.008\pm0.001	94.17 \pm 0.15	0.214 \pm 0.004	0.031 \pm 0.001	0.179 \pm 0.003	0.007 \pm 0.001
+ TDiv-SDiv	93.95 \pm 0.01	0.205\pm0.006	0.028\pm0.001	0.181\pm0.005	0.008\pm0.002	94.19\pm0.11	0.200\pm0.004	0.027\pm0.002	0.178\pm0.002	0.006\pm0.002
<i>Multiple forward passes for inference:</i>										
KD + BE (Mariet et al., 2021)	93.77 \pm 0.20	0.275 \pm 0.011	0.043 \pm 0.002	0.201 \pm 0.006	0.011 \pm 0.001	94.00 \pm 0.05	0.263 \pm 0.008	0.040 \pm 0.001	0.195 \pm 0.003	0.010 \pm 0.001
+ ConfODS (Nam et al., 2021)	94.06 \pm 0.13	0.223 \pm 0.007	0.031 \pm 0.002	0.188 \pm 0.004	0.006 \pm 0.001	94.06 \pm 0.10	0.222 \pm 0.003	0.032 \pm 0.001	0.185 \pm 0.002	0.008 \pm 0.001
+ TDiv-SDiv	93.74 \pm 0.10	0.213 \pm 0.002	0.028 \pm 0.001	0.189 \pm 0.001	0.006 \pm 0.001	94.10 \pm 0.13	0.202 \pm 0.002	0.026 \pm 0.000	0.181 \pm 0.002	0.007 \pm 0.001

Table 3. Results of the distilled students for WRN28x4 on CIFAR-100. Results with \pm std. are averaged over 4 seeds.

Method	Students distilled from DE-2 teacher					Students distilled from DE-4 teacher				
	ACC (\uparrow)	NLL (\downarrow)	ECE (\downarrow)	cNLL (\downarrow)	cECE (\downarrow)	ACC (\uparrow)	NLL (\downarrow)	ECE (\downarrow)	cNLL (\downarrow)	cECE (\downarrow)
<i>Single forward pass for inference:</i>										
KD (Hinton et al., 2015)	79.15 \pm 0.25	1.029 \pm 0.008	0.125 \pm 0.002	0.864 \pm 0.006	0.046 \pm 0.004	79.09 \pm 0.20	1.038 \pm 0.011	0.130 \pm 0.003	0.861 \pm 0.007	0.046 \pm 0.001
AE-KD (Du et al., 2020)	78.79 \pm 0.23	1.041 \pm 0.015	0.129 \pm 0.002	0.871 \pm 0.009	0.044 \pm 0.004	79.00 \pm 0.39	1.033 \pm 0.009	0.129 \pm 0.004	0.859 \pm 0.008	0.045 \pm 0.004
Proxy-EnD ² (Ryabinin et al., 2021)	78.40 \pm 0.28	1.072 \pm 0.012	0.138 \pm 0.003	0.894 \pm 0.007	0.047 \pm 0.002	78.75 \pm 0.28	1.076 \pm 0.016	0.138 \pm 0.002	0.886 \pm 0.011	0.046 \pm 0.003
KD + LatentBE (Ours)	79.16\pm0.18	0.985\pm0.011	0.122\pm0.002	0.848\pm0.006	0.045\pm0.003	79.46\pm0.20	0.993\pm0.024	0.124\pm0.004	0.837\pm0.012	0.046\pm0.005
+ ConfODS	78.61 \pm 0.19	0.965 \pm 0.008	0.109 \pm 0.001	0.873 \pm 0.006	0.046 \pm 0.002	79.27 \pm 0.30	0.955 \pm 0.008	0.115 \pm 0.002	0.840 \pm 0.007	0.048 \pm 0.004
+ TDiv-SDiv	79.49\pm0.15	0.826\pm0.007	0.072\pm0.003	0.798\pm0.005	0.041\pm0.002	80.02\pm0.07	0.792\pm0.004	0.067\pm0.001	0.772\pm0.003	0.041\pm0.003
<i>Multiple forward passes for inference:</i>										
KD + BE (Mariet et al., 2021)	78.50 \pm 0.42	1.067 \pm 0.010	0.134 \pm 0.003	0.888 \pm 0.008	0.044 \pm 0.003	78.92 \pm 0.24	1.035 \pm 0.013	0.130 \pm 0.002	0.863 \pm 0.012	0.043 \pm 0.002
+ ConfODS (Nam et al., 2021)	78.24 \pm 0.19	1.011 \pm 0.013	0.118 \pm 0.002	0.897 \pm 0.005	0.048 \pm 0.003	78.65 \pm 0.29	1.002 \pm 0.013	0.123 \pm 0.002	0.873 \pm 0.010	0.046 \pm 0.003
+ TDiv-SDiv	78.50 \pm 0.21	0.871 \pm 0.005	0.080 \pm 0.003	0.837 \pm 0.004	0.044 \pm 0.002	79.56 \pm 0.18	0.818 \pm 0.007	0.066 \pm 0.002	0.798 \pm 0.006	0.042 \pm 0.003

5.3. Results on CIFAR-10/100

We compare ours to the existing ensemble distillation methods. Here, we consider baselines using a single student network: KD (Hinton et al., 2015), AE-KD (Du et al., 2020), and Proxy-EnD² (Ryabinin et al., 2021). As suggested in Ashukha et al. (2020), we report both original and calibrated metrics for NLL and ECE. See Appendix B.4 for the details in distillation methods, and Appendix B.3 for the definitions of evaluation metrics.

The results for CIFAR-10 and CIFAR-100 are presented in Tables 2 and 3. LatentBE consistently outperforms the baselines both in terms of accuracy and uncertainty estimates, and our perturbation further boosts up the performance of LatentBE. We note that our method gets better as the number of teachers increases, indicating that ours effectively transfers diversities from multiple teachers.

Moreover, a benefit of our approach is that it consistently outperforms the vanilla KD even when the number of teachers M is small. On the other hand, Ryabinin et al. (2021) suffers when M is small because the estimation for Dirichlet parameters required for the distillation become inaccurate. This gives an advantage to our method under resource limited setting where training large number of ensemble teachers are intractable.

Table 4. Results of the distilled students for WRN28x1 on CIFAR-10-C, and WRN28x4 on CIFAR-100-C. The results are with the DE-4 teacher and are averaged over 4 seeds.

Method	ACC (\uparrow)	NLL (\downarrow)	ECE (\downarrow)
CIFAR-10-C:			
<u>DE-4 teacher</u>	<u>73.18</u>	<u>1.025</u>	<u>0.092</u>
KD (Hinton et al., 2015)	72.49 \pm 0.38	1.492 \pm 0.039	0.202 \pm 0.004
AE-KD (Du et al., 2020)	72.06 \pm 0.49	1.540 \pm 0.046	0.207 \pm 0.005
Proxy-EnD ² (Ryabinin et al., 2021)	71.05 \pm 0.51	1.600 \pm 0.050	0.217 \pm 0.007
KD + LatentBE (Ours)	72.73\pm0.54	1.522\pm0.078	0.203\pm0.008
+ ConfODS	70.18 \pm 0.48	1.477 \pm 0.038	0.200 \pm 0.004
+ TDiv-SDiv	73.22\pm0.43	1.237\pm0.024	0.171\pm0.004
CIFAR-100-C:			
<u>DE-4 teacher</u>	<u>51.08</u>	<u>2.296</u>	<u>0.114</u>
KD (Hinton et al., 2015)	48.79 \pm 0.16	3.410 \pm 0.056	0.341 \pm 0.005
AE-KD (Du et al., 2020)	48.86 \pm 0.12	3.407 \pm 0.066	0.340 \pm 0.007
Proxy-EnD ² (Ryabinin et al., 2021)	48.67 \pm 0.25	3.378 \pm 0.054	0.350 \pm 0.004
KD + LatentBE (Ours)	50.03\pm0.36	3.178\pm0.053	0.321\pm0.003
+ ConfODS	47.17 \pm 0.12	3.273 \pm 0.019	0.326 \pm 0.002
+ TDiv-SDiv	51.36\pm0.28	2.507\pm0.059	0.208\pm0.006

Robustness to common corruptions We further evaluate our method on corrupted CIFAR datasets to verify its robustness under common corruptions (Hendrycks & Dietterich, 2019). Table 4 reports the evaluation metrics averaged over all types of corruptions and intensities and shows that ours are better calibrated than the existing baselines.

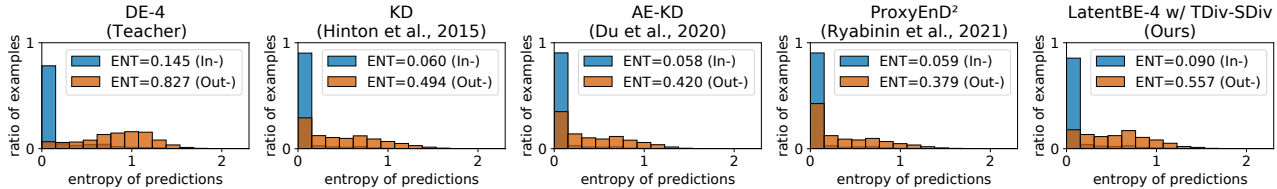


Figure 5. Histograms of the predictive entropy for WRN28x1 students distilled from DE-4 on CIFAR-10. ‘In-’ denotes the entropy on the in-distribution data (i.e., CIFAR-10), and ‘Out-’ denotes the entropy on the out-of-distribution data (i.e., SVHN).

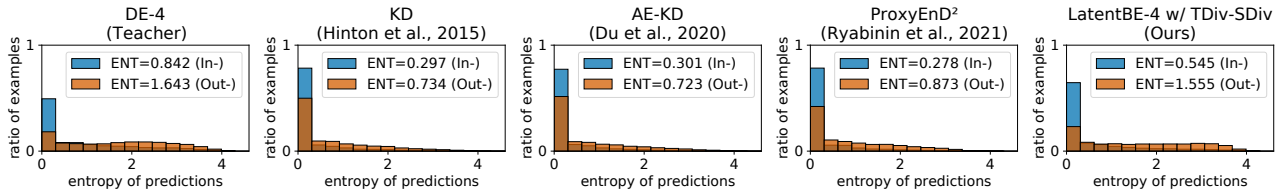


Figure 6. Histograms of the predictive entropy for WRN28x4 students distilled from DE-4 on CIFAR-100. ‘In-’ denotes the entropy on the in-distribution data (i.e., CIFAR-100), and ‘Out-’ denotes the entropy on the out-of-distribution data (i.e., SVHN).

Out-of-distribution data A well-calibrated classifier should be uncertain for out-of-distribution data while being certain for in-distribution data. Following Lakshminarayanan et al. (2017), we draw histograms depicting the distribution of the predictive entropy for in-distribution (i.e., CIFAR-10/100) and out-of-distribution data (i.e., SVHN) in Figures 5 and 6. It shows that ours are more uncertain about unseen classes, behaving similarly to ensemble teachers. Refer to Appendix A.2 for more results with other out-of-distribution datasets.

5.4. Runtime analysis

Tables 2 and 3 show that LatentBE is competitive or even better for some metrics than the one-to-one distillation methods with BE students (Mariet et al., 2021; Nam et al., 2021). This is quite remarkable since BE requires additional inference cost while LatentBE doesn’t. More specifically, Table 5 reports the runtimes of BE and LatentBE on the same single GeForce RTX 3090 setting. We measure the wall-clock time in hours for training and milliseconds per image (ms/img) for inference. In principle, BE requires multiple forward passes for inference (i.e., 0.390 ms/img), but this can be avoided by parallelized inference (0.288 ms/img). Still, the parallelization requires additional cost from handling larger mini-batches, which is still significantly larger than a single model inference time (0.107 ms/img).

5.5. Results on TinyImageNet and ImageNet-1k

We verify the scalability of our proposed method on large-scale datasets, including TinyImageNet and ImageNet-1k (Russakovsky et al., 2015). Tables 6 and 7 show that ours outperform the baselines for both datasets.

Table 5. Comparison on training and testing runtime. The results are with WRN28x4 on CIFAR-100 previously reported in Table 3.

Method	Performance		Runtime	
	ACC (↑)	NLL (↓)	Training (↓)	Inference (↓)
DE-4 teacher	81.37	0.706	6.2 hrs.	0.390 ms/img
KD (Hinton et al., 2015)	79.09	1.038	1.6 hrs.	0.107 ms/img
KD + LatentBE (Ours)	79.46	0.993	3.3 hrs.	0.107 ms/img
+ TDiv-SDiv	80.02	0.792	5.3 hrs.	0.107 ms/img
KD + BE (Mariet et al., 2021)	78.92	1.035	3.3 hrs.	0.288 ms/img
+ TDiv-SDiv	79.56	0.818	5.3 hrs.	0.288 ms/img

6. Conclusion

In this paper, we proposed a novel ensemble distillation algorithm improving both prediction accuracy and uncertainty calibration without increasing inference cost. We first presented LatentBE, where the rank-one factors of BEs are trained in a one-to-one way, but later weight-averaged for inference. We showed that under a suitable training scheme, the subspaces defined by the rank-one factors of BE remain in a flat minimum, and weight averaging those rank-one factors thus yields a robust single student model. We further presented a novel perturbation strategy for ensemble distillation that decreases student diversities and increases teacher diversities at the same time. By training with inputs perturbed in that way, we can effectively enhance the diversities of students. Our ensemble distillation algorithm combining these two achieved remarkable performance on various image classification tasks.

Acknowledgements

This work was partly supported by KAIST-NAVER Hypercreative AI Center, Institute of Information & communications Technology Planning & Evaluation (IITP) grant

Table 6. Results of R18 and WRN28x4 on TinyImageNet. Results with \pm std. are averaged over 3 seeds.

Method	R18 on TinyImageNet					WRN28x4 on TinyImageNet				
	ACC (\uparrow)	NLL (\downarrow)	ECE (\downarrow)	cNLL (\downarrow)	cECE (\downarrow)	ACC (\uparrow)	NLL (\downarrow)	ECE (\downarrow)	cNLL (\downarrow)	cECE (\downarrow)
Baseline results:										
Base (w/o distillation)	64.78 \pm 0.12	1.580 \pm 0.004	0.104 \pm 0.002	1.494 \pm 0.001	0.031 \pm 0.002	63.25 \pm 0.07	1.599 \pm 0.010	0.099 \pm 0.003	1.517 \pm 0.006	0.022 \pm 0.001
Single KD (Hinton et al., 2015)	67.29 \pm 0.16	1.470 \pm 0.005	0.107 \pm 0.002	1.383 \pm 0.003	0.046 \pm 0.002	66.25 \pm 0.38	1.443 \pm 0.012	0.077 \pm 0.005	1.391 \pm 0.007	0.027 \pm 0.002
DE-4 teacher	69.28	1.273	0.025	1.272	0.023	68.75	1.276	0.027	1.277	0.027
Students distilled from DE-4 teacher:										
KD (Hinton et al., 2015)	68.88 \pm 0.20	1.391 \pm 0.009	0.099 \pm 0.001	1.317 \pm 0.006	0.044 \pm 0.002	67.69 \pm 0.08	1.351 \pm 0.004	0.067 \pm 0.002	1.314 \pm 0.004	0.026 \pm 0.006
AE-KD (Du et al., 2020)	67.20 \pm 0.08	1.409 \pm 0.006	0.089 \pm 0.001	1.355 \pm 0.005	0.040 \pm 0.003	64.77 \pm 0.08	1.424 \pm 0.004	0.052 \pm 0.002	1.403 \pm 0.004	0.020 \pm 0.006
Proxy-EnD ² (Ryabinin et al., 2021)	62.42 \pm 0.26	1.572 \pm 0.002	0.017 \pm 0.003	1.571 \pm 0.002	0.021 \pm 0.004	62.29 \pm 0.25	1.578 \pm 0.003	0.049 \pm 0.004	1.560 \pm 0.005	0.016 \pm 0.003
KD + LatentBE (Ours)	68.96 \pm 0.19	1.391 \pm 0.007	0.103 \pm 0.002	1.312 \pm 0.007	0.042 \pm 0.002	67.76 \pm 0.05	1.343 \pm 0.004	0.073 \pm 0.001	1.303 \pm 0.004	0.025 \pm 0.001
+ ConfODS	69.00 \pm 0.32	1.390 \pm 0.008	0.101 \pm 0.004	1.313 \pm 0.005	0.047 \pm 0.005	67.89 \pm 0.18	1.338 \pm 0.005	0.071 \pm 0.003	1.299 \pm 0.004	0.027 \pm 0.002
+ TDiv-SDiv	69.14 \pm 0.27	1.342 \pm 0.005	0.085 \pm 0.004	1.290 \pm 0.001	0.038 \pm 0.002	68.15 \pm 0.03	1.317 \pm 0.008	0.062 \pm 0.000	1.286 \pm 0.006	0.027 \pm 0.002

Table 7. Results for R50 on ImageNet-1k.

Method	R50 on ImageNet-1k				
	ACC (\uparrow)	NLL (\downarrow)	ECE (\downarrow)	cNLL (\downarrow)	cECE (\downarrow)
Baselines results:					
Base (w/o distillation)	76.80	0.927	0.040	0.913	0.019
Single KD (Hinton et al., 2015)	76.90	0.918	0.028	0.913	0.017
DE-2 teacher	77.96	0.862	0.018	0.859	0.018
Students distilled from DE-2 teacher:					
KD (Hinton et al., 2015)	77.01	0.904	0.028	0.900	0.017
KD + LatentBE (Ours)	77.30	0.902	0.029	0.895	0.019
+ TDiv-SDiv	77.38	0.898	0.027	0.892	0.018

funded by the Korea government (MSIT) (No.2019-0-00075, Artificial Intelligence Graduate School Program (KAIST), No. 2021-0-02068, Artificial Intelligence Innovation Hub), and National Research Foundation of Korea (NRF) funded by the Ministry of Education (NRF-2021M3E5D9025030).

References

Ashukha, A., Lyzhov, A., Molchanov, D., and Vetrov, D. P. Pitfalls of in-domain uncertainty estimation and ensembling in deep learning. In *International Conference on Learning Representations (ICLR)*, 2020.

D’Angelo, F. and Fortuin, V. Repulsive deep ensembles are bayesian. In *Advances in Neural Information Processing Systems 34 (NeurIPS 2021)*, 2021.

Du, S., You, S., Li, X., Wu, J., Wang, F., Qian, C., and Zhang, C. Agree to disagree: Adaptive ensemble knowledge distillation in gradient space. In *Advances in Neural Information Processing Systems 33 (NeurIPS 2020)*, 2020.

Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. On calibration of modern neural networks. In *Proceedings of The 34th International Conference on Machine Learning (ICML 2017)*, 2017.

He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE*

Conference on Computer Vision and Pattern Recognition (CVPR), 2016.

Hendrycks, D. and Dietterich, T. Benchmarking neural network robustness to common corruptions and perturbations. *International Conference on Learning Representations (ICLR)*, 2019.

Hinton, G. E., Vinyals, O., and Dean, J. Distilling the knowledge in a neural network. In *Deep Learning and Representation Learning Workshop, NIPS 2014*, 2015.

Izmailov, P., Podoprikin, D., Garipov, T., Vetrov, D., and Wilson, A. G. Averaging weights leads to wider optima and better generalization. In *Proceedings of the 34th Conference on Uncertainty in Artificial Intelligence (UAI 2018)*, 2018.

Jin, X., Peng, B., Wu, Y., Liu, Y., Liu, J., Liang, D., Yan, J., and Hu, X. Knowledge distillation via route constrained optimization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019.

Krizhevsky, A. Learning multiple layers of features from tiny images. In *Citeseer*, 2009.

Krogh, A. and Hertz, J. A simple weight decay can improve generalization. In *Advances in Neural Information Processing Systems 4 (NIPS 1991)*, 1991.

Lakshminarayanan, B., Pritzel, A., and Blundell, C. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems 30 (NIPS 2017)*, 2017.

Malinin, A., Mlodozieniec, B., and Gales, M. Ensemble distribution distillation. In *International Conference on Learning Representations (ICLR)*, 2020.

Mariet, Z. E., Jenatton, R., Wenzel, F., and Tran, D. Distilling ensembles improves uncertainty estimates. In *Third Symposium on Advances in Approximate Bayesian Inference*, 2021.

- Mirzadeh, S. I., Farajtabar, M., Li, A., Levine, N., Matsukawa, A., and Ghasemzadeh, H. Improved knowledge distillation via teacher assistant. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- Nam, G., Yoon, J., Lee, Y., and Lee, J. Diversity matters when learning from ensembles. In *Advances in Neural Information Processing Systems 34 (NeurIPS 2021)*, 2021.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011.
- Ovadia, Y., Fertig, E., Ren, J., Nado, Z., Sculley, D., Nowozin, S., Dillon, J. V., Lakshminarayanan, B., and Snoek, J. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, 2019.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)*, 2019.
- Rame, A. and Cord, M. Dice: Diversity in deep ensembles via conditional redundancy adversarial estimation. In *International Conference on Learning Representations (ICLR)*, 2021.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., and Fei-Fei, L. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision (IJCV)*, 2015.
- Ryabinin, M., Malinin, A., and Gales, M. Scaling ensemble distribution distillation to many classes with proxy targets. In *Advances in Neural Information Processing Systems 34 (NeurIPS 2021)*, 2021.
- Srinivas, S. and Fleuret, F. Knowledge transfer with jacobian matching. In *Proceedings of The 35th International Conference on Machine Learning (ICML 2018)*, 2018.
- Tashiro, Y., Song, Y., and Ermon, S. Diversity can be transferred: Output diversification for white- and black-box attacks. In *Advances in Neural Information Processing Systems 33 (NeurIPS 2020)*, 2020.
- Tran, L., Veeling, B. S., Roth, K., Swiatkowski, J., Dillon, J. V., Snoek, J., Mandt, S., Salimans, T., Nowozin, S., and Jenatton, R. Hydra: Preserving ensemble diversity for model distillation. *arXiv:2001.04694*, 2020.
- Wen, Y., Tran, D., and Ba, J. Batchensemble: an alternative approach to efficient ensemble and lifelong learning. In *International Conference on Learning Representations (ICLR)*, 2020.
- Wenzel, F., Snoek, J., Tran, D., and Jenatton, R. Hyperparameter ensembles for robustness and uncertainty quantification. In *Advances in Neural Information Processing Systems 33 (NeurIPS 2020)*, 2020.
- Wortsman, M., Horton, M., Guestrin, C., Farhadi, A., and Rastegari, M. Learning neural network subspaces. In *Advances in Neural Information Processing Systems 34 (NeurIPS 2021)*, 2021.
- Yu, F., Seff, A., Zhang, Y., Song, S., Funkhouser, T., and Xiao, J. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.
- Zagoruyko, S. and Komodakis, N. Wide residual networks. In *Proceedings of the British Machine Vision Conference (BMVC)*, 2016.
- Zaidi, S., Zela, A., Elsken, T., Holmes, C. C., Hutter, F., and Teh, Y. Neural ensemble search for uncertainty estimation and dataset shift. In *Advances in Neural Information Processing Systems 34 (NeurIPS 2021)*, 2021.

Table 8. Baseline results for the experiments on CIFAR-10/100. Results with \pm std. are averaged over 4 seeds.

Method	WRN28x1 on CIFAR-10					WRN28x4 on CIFAR-100				
	ACC (\uparrow)	NLL (\downarrow)	ECE (\downarrow)	cNLL (\downarrow)	cECE (\downarrow)	ACC (\uparrow)	NLL (\downarrow)	ECE (\downarrow)	cNLL (\downarrow)	cECE (\downarrow)
DE-1	93.01	0.271	0.038	0.222	0.009	78.03	0.888	0.062	0.878	0.041
DE-2	94.11	0.203	0.016	0.191	0.010	80.13	0.768	0.026	0.768	0.025
DE-3	94.53	0.180	0.011	0.176	0.012	80.93	0.723	0.026	0.720	0.022
DE-4	94.71	0.170	0.006	0.168	0.010	81.37	0.706	0.031	0.700	0.018
DE-5	94.68	0.166	0.008	0.165	0.009	-	-	-	-	-
DE-6	94.64	0.162	0.007	0.162	0.007	-	-	-	-	-
DE-7	94.89	0.159	0.008	0.159	0.008	-	-	-	-	-
DE-8	95.07	0.157	0.009	0.157	0.009	-	-	-	-	-
Base (w/o distillation)	93.05 \pm 0.16	0.263 \pm 0.008	0.037 \pm 0.002	0.218 \pm 0.005	0.006 \pm 0.002	77.93 \pm 0.21	0.893 \pm 0.006	0.060 \pm 0.002	0.882 \pm 0.006	0.041 \pm 0.003
Single KD (Hinton et al., 2015)	93.52 \pm 0.09	0.274 \pm 0.008	0.041 \pm 0.001	0.206 \pm 0.004	0.009 \pm 0.001	78.47 \pm 0.13	1.039 \pm 0.005	0.127 \pm 0.003	0.886 \pm 0.003	0.045 \pm 0.002

A. Additional Results

A.1. Baseline results for CIFAR-10/100

In Table 8, we report the baseline results for the experiments on CIFAR-10/100: (1) evaluation results for DE teachers distilling knowledge into students, (2) performance of the single model trained with the classical cross-entropy loss without distillation, and (3) performance of the single model distilled from DE-1.

A.2. Further experiments on predictive uncertainty

For the experiments on predictive uncertainty, we consider SVHN (Netzer et al., 2011), LSUN (Yu et al., 2015), and TinyImageNet (Russakovsky et al., 2015) as out-of-distribution data. Here, LSUN and TinyImageNet images are downscaled into $32 \times 32 \times 3$. Figures 7 and 8 further provide the predictive uncertainty results on LSUN and TinyImageNet. Again, our approach exhibits higher predictive uncertainty on out-of-distribution examples than existing baselines.

B. Experimental Details

Code is available at <https://github.com/cs-giung/distill-latentbe>. Our implementation for the experiments on CIFAR-10/100 and TinyImageNet are built on PyTorch (Paszke et al., 2019). Besides, the experiments on ImageNet-1k are conducted with 8 TPUv3 cores, supported by the TPU Research Cloud².

B.1. Datasets

CIFAR-10/100 The dataset is available at <https://www.cs.toronto.edu/~kriz/cifar.html>. It consists of 50,000 train examples and 10,000 test examples from 10/100 classes, with images size of $32 \times 32 \times 3$. In this paper, the last 5,000 examples of the train split are used as the validation split for computing calibrated metrics. We follow the standard data augmentation policy (He et al., 2016) which consists of random cropping of 32 pixels with a padding of 4 pixels and random horizontal flipping. Throughout experiments on CIFAR-10/100 classification, we use WideResNet (WRN) networks introduced in Zagoruyko & Komodakis (2016); WRN28x1 on CIFAR-10 and WRN28x4 on CIFAR-100.

TinyImageNet The dataset is available at <http://cs231n.stanford.edu/tiny-imagenet-200.zip>. It consists of 100,000 train examples, 10,000 validation examples and 10,000 test examples from 200 classes subsampled from ImageNet-1k, with images size of $64 \times 64 \times 3$. Since the labels of the official test set are not publicly available, we use the official validation set as a test set for experiments. Consequently, the last 500 examples for each class of the train split are used as the validation split for computing calibrated metrics, i.e., train and validation split consists of 90,000 and 10,000 examples, respectively. We apply the data augmentation which consists of random cropping of 64 pixels with a padding of 4 pixels and random horizontal flipping. Throughout experiments on TinyImageNet classification, we use the ResNet-18 (R18) network introduced in He et al. (2016) and the WRN28x4 network introduced in Zagoruyko & Komodakis (2016).

²<https://sites.research.google/trc/about/>

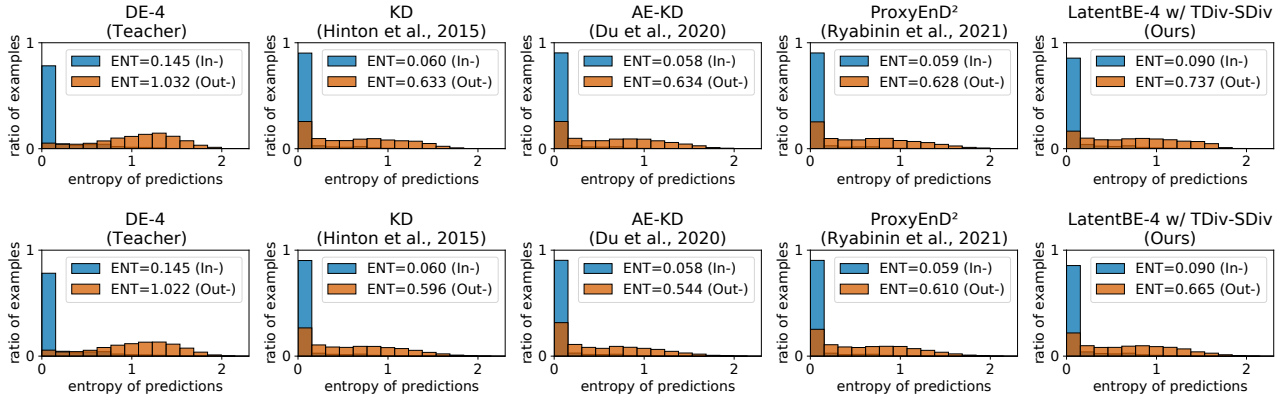


Figure 7. Histograms of the predictive entropy for WRN28x1 students distilled from DE-4 on CIFAR-10. ‘In-’ denotes the entropy on the in-distribution data (i.e., CIFAR-10), and ‘Out-’ denotes the entropy on the out-of-distribution data including LSUN (1st row) and TinyImageNet (2nd row).

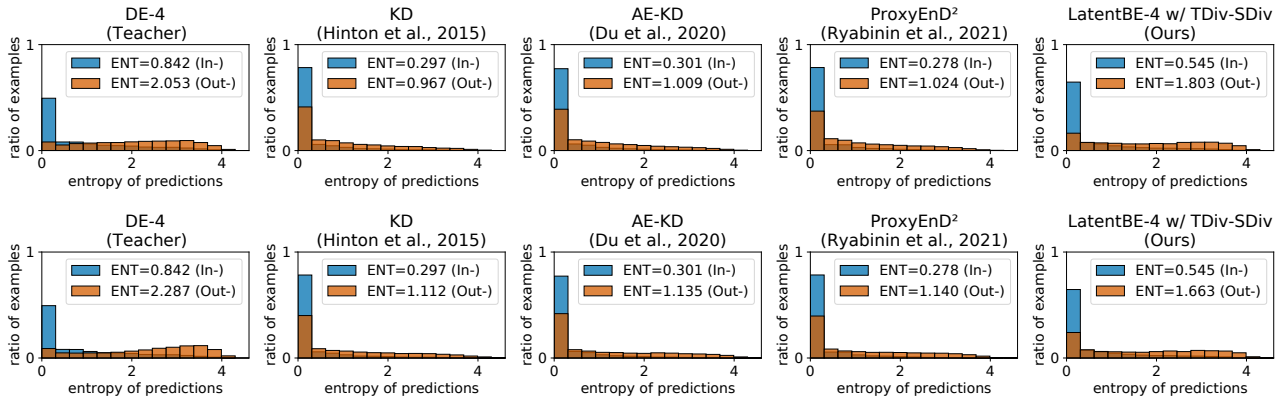


Figure 8. Histograms of the predictive entropy for WRN28x4 students distilled from DE-4 on CIFAR-100. ‘In-’ denotes the entropy on the in-distribution data (i.e., CIFAR-100), and ‘Out-’ denotes the entropy on the out-of-distribution data including LSUN (1st row) and TinyImageNet (2nd row).

ImageNet-1k It consists of 1,281,167 train examples, 50,000 validation examples and 100,000 test images from 1,000 classes. Since the labels of the official test set are not publicly available, we only report the evaluation results on the validation set. We follow the standard data augmentation policy from PyTorch which consists of random cropping with an images size of $224 \times 224 \times 3$ and random horizontal flipping³. Throughout experiments on ImageNet-1k classification, we use the ResNet-50 (R50) network introduced in He et al. (2016).

B.2. Training details

All images are standardized by subtracting the per-channel mean and dividing the result by the per-channel standard deviation. We use SGD optimizer with Nesterov momentum 0.9, and a single-cycle cosine annealing learning rate schedule with a linear warm-up, i.e., the learning rate starts from $0.01 \times \text{base_lr}$ and reaches base_lr after the first 5 epochs, and is decayed by the single-cycle cosine annealing learning rate schedule. More precisely, (1) for CIFAR-10/100, we run 200 epochs on a single machine with batch size 128 and $\text{base_lr} = 0.1$, (2) for TinyImageNet, we run 80 epochs on four machines with the total batch size 128 and $\text{base_lr} = 0.1$, and (3) for ImageNet-1k, we run 100 epochs on eight machines with the total batch size 256 and $\text{base_lr} = 0.1$. We also apply the weight decay (Krogh & Hertz, 1991) to regularize training; the weight decay coefficient is set to be 0.0005 for CIFAR-10/100 and TinyImageNet, and 0.0001 for ImageNet-1k.

³<https://github.com/pytorch/examples/tree/master/imagenet>

B.3. Evaluation

The problem we address in this paper is the K -way classification problem; a neural network $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^K$ takes D -dimensional inputs \mathbf{x} (i.e., images) and makes predictions about outputs y (i.e., class label) with K -dimensional logits. We denote the output probabilities of the model \mathbf{f} for a given input \mathbf{x} as

$$\mathbf{p}_{\mathbf{f}}^{(k)}(\mathbf{x}) = \frac{\exp(\mathbf{f}^{(k)}(\mathbf{x}))}{\sum_{j=1}^K \exp(\mathbf{f}^{(j)}(\mathbf{x}))}, \quad \text{for } k = 1, \dots, K. \quad (18)$$

Standard metrics We evaluate the following *standard metrics* of the model \mathbf{f} on the dataset \mathcal{D} :

- ACC (accuracy; higher is better):

$$\text{ACC}(\mathbf{f}, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}} \left[y = \arg \max_k \mathbf{p}_{\mathbf{f}}^{(k)}(\mathbf{x}) \right], \quad (19)$$

where $[\cdot]$ denotes the Iverson bracket.

- NLL (negative log-likelihood; lower is better):

$$\text{NLL}(\mathbf{f}, \mathcal{D}) = - \sum_{(\mathbf{x}, y) \in \mathcal{D}} y \log \mathbf{p}_{\mathbf{f}}^{(y)}(\mathbf{x}) \quad (20)$$

- ECE (expected calibration error; lower is better):

$$\text{ECE}(\mathbf{f}, \mathcal{D}, L) = \sum_{l=1}^L \frac{|\mathcal{B}_l|}{|\mathcal{D}|} \left| \text{ACC}(\mathbf{f}, \mathcal{B}_l) - \sum_{(\mathbf{x}, \cdot) \in \mathcal{B}_l} \frac{\max_k \mathbf{p}_{\mathbf{f}}^{(k)}(\mathbf{x})}{|\mathcal{B}_l|} \right|, \quad (21)$$

where $\{\mathcal{B}_1, \dots, \mathcal{B}_L\}$ is a partition of \mathcal{D} , where $\mathcal{B}_l = \{(\mathbf{x}, y) \in \mathcal{D} \mid \max_k \mathbf{p}_{\mathbf{f}}^{(k)}(\mathbf{x}) \in ((l-1)/L, l/L]\}$. Here, the difference between the accuracy and mean confidence of predictions represents the calibration gap for each bin. We fixed $L = 15$ for all evaluation results in this paper.

Calibrated metrics Temperature scaling softens output probabilities of the model with a single scale parameter $\tau > 0$, and it can be used for calibrating probabilistic models *without affecting the model's prediction accuracy* (Guo et al., 2017). We define the temperature scaled output probabilities of the model \mathbf{f} for a given input \mathbf{x} as

$$\mathbf{p}_{\mathbf{f}}^{(k)}(\mathbf{x}; \tau) = \frac{\exp(\mathbf{f}^{(k)}(\mathbf{x})/\tau)}{\sum_{j=1}^K \exp(\mathbf{f}^{(j)}(\mathbf{x})/\tau)}, \quad \text{for } k = 1, \dots, K. \quad (22)$$

Following Ashukha et al. (2020), we also evaluate the *calibrated metrics* which are computed using the temperature scaled outputs. Specifically, we first find the optimal temperature which minimizes the NLL on the validation split $\mathcal{D}_{\text{valid}}$, i.e.,

$$\tau^* \leftarrow \arg \min_{\tau} \left[- \sum_{(\mathbf{x}, y) \in \mathcal{D}_{\text{valid}}} y \log \mathbf{p}_{\mathbf{f}}^{(y)}(\mathbf{x}; \tau) \right], \quad (23)$$

and then we can compute the following *calibrated metrics*:

- cNLL (calibrated negative log-likelihood):

$$\text{cNLL}(\mathbf{f}, \mathcal{D}, \tau^*) = - \sum_{(\mathbf{x}, y) \in \mathcal{D}} y \log \mathbf{p}_{\mathbf{f}}^{(y)}(\mathbf{x}; \tau^*) \quad (24)$$

- cECE (calibrated expected calibration error):

$$\text{cECE}(\mathbf{f}, \mathcal{D}, L, \tau^*) = \sum_{l=1}^L \frac{|\mathcal{B}_l|}{|\mathcal{D}|} \left| \text{ACC}(\mathbf{f}, \mathcal{B}_l) - \sum_{(\mathbf{x}, \cdot) \in \mathcal{B}_l} \frac{\max_k \mathbf{p}_{\mathbf{f}}^{(k)}(\mathbf{x}; \tau^*)}{|\mathcal{B}_l|} \right|. \quad (25)$$

B.4. Ensemble Distillation Methods

Ensemble distillation (KD) Let $\{\mathcal{T}_1, \dots, \mathcal{T}_M\}$ be a set of pre-trained teachers, and \mathcal{S}_θ be a student. In practice, we often consider the cross-entropy loss to ground-truth labels during training in addition to the KD loss defined in Equation (2). For a given image \mathbf{x} and a corresponding one-hot class label \mathbf{y} , the loss for the KD with the cross-entropy term is defined as

$$(1 - \alpha)\mathcal{H}[\mathbf{y}, \mathbf{p}_{\mathcal{S}_\theta}(\mathbf{x}; \tau)] + \alpha\tau^2 \frac{1}{M} \sum_{m=1}^M \mathcal{H}[\mathbf{p}_{\mathcal{T}_m}(\mathbf{x}; \tau), \mathbf{p}_{\mathcal{S}_\theta}(\mathbf{x}; \tau)]. \quad (26)$$

Here, we have *two hyperparameters*: (1) τ smooth output probabilities via temperature scaling, and (2) α adjusts the balance between two cross-entropy losses. We use $(\alpha, \tau) = (1.0, 4.0)$ for experiments on CIFAR-10/100, $(\alpha, \tau) = (0.9, 20.0)$ for experiments on TinyImageNet, and $(\alpha, \tau) = (1.0, 1.0)$ for experiments on ImageNet-1k.

Adaptive ensemble knowledge distillation (AE-KD) Du et al. (2020) argued that the vanilla ensemble distillation with Equation (2) produces the learning signal determined by the most dominant teacher. To resolve the issue, they propose Adaptive Ensemble Knowledge Distillation (AE-KD) that optimizes weighting coefficients $\{\omega_m\}_{m=1}^M$ of the weighted ensemble distillation loss,

$$\sum_{m=1}^M \omega_m \mathcal{H}[\mathbf{p}_{\mathcal{T}_m}(\mathbf{x}; \tau), \mathbf{p}_{\mathcal{S}_\theta}(\mathbf{x}; \tau)], \quad (27)$$

by solving the following optimization problem:

$$\min_{\omega_1, \dots, \omega_M} \frac{1}{2\tau^2} \left\| \mathbf{p}_{\mathcal{S}_\theta}(\mathbf{x}) - \sum_{m=1}^M \omega_m \mathbf{p}_{\mathcal{T}_m}(\mathbf{x}) \right\|_2^2 \quad (28)$$

$$\text{subject to } \sum_{m=1}^M \omega_m = 1, \quad 0 \leq \omega_m \leq C \text{ for } m = 1, \dots, M. \quad (29)$$

It introduces an *additional hyperparameter*: $C \in [1/M, 1]$ controls the *tolerance of disagreement* among teachers, i.e., with the decrease of C , more tolerance of disagreement among the gradients is allowed. Note that the vanilla ensemble distillation is a specialization of AE-KD when $C = 1/M$. Throughout all experiments, we use $C = 0.6$.

Ensemble distribution distillation with Proxy-Dirichlet distribution (Proxy-EnD²) Ryabinin et al. (2021) introduce a Proxy-Dirichlet target having the density of

$$\mathbf{q}_{\mathcal{T}}(\mathbf{p}|\mathbf{x}) = \frac{\Gamma\left(\sum_{j=1}^K \beta_j(\mathbf{x})\right)}{\prod_{i=1}^K \Gamma(\beta_i(\mathbf{x}))} \prod_{k=1}^K \left(\mathbf{p}^{(k)}\right)^{\beta_k(\mathbf{x})-1}, \quad (30)$$

where the concentration parameters are approximated from output probabilities of the teachers,

$$\beta_k(\mathbf{x}) \leftarrow \mathbf{p}_{\mathcal{T}}^{(k)}(\mathbf{x}) \frac{(K-1)/2}{\sum_{j=1}^K \left[\mathbf{p}_{\mathcal{T}}^{(j)}(\mathbf{x}) \left(\log \mathbf{p}_{\mathcal{T}}^{(j)}(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^M \log \mathbf{p}_{\mathcal{T}_m}^{(j)}(\mathbf{x}) \right) \right]}. \quad (31)$$

Here, $\mathbf{p}_{\mathcal{T}}$ denotes the mean prediction of the teachers, i.e., $\mathbf{p}_{\mathcal{T}}^{(k)}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M \mathbf{p}_{\mathcal{T}_m}^{(k)}(\mathbf{x})$. This *ensemble distribution* is distilled into the student network \mathcal{S}_θ which represents the density over $(K-1)$ -simplex,

$$\mathbf{q}_{\mathcal{S}_\theta}(\mathbf{p}|\mathbf{x}) = \frac{\Gamma\left(\sum_{j=1}^K \mathcal{S}_\theta^{(j)}(\mathbf{x})\right)}{\prod_{i=1}^K \Gamma(\mathcal{S}_\theta^{(i)}(\mathbf{x}))} \prod_{k=1}^K \left(\mathbf{p}^{(k)}\right)^{\mathcal{S}_\theta^{(k)}(\mathbf{x})-1}, \quad (32)$$

where the student network represents the concentration parameters to model the Dirichlet distribution, i.e., the k^{th} concentration parameter is $e^{\mathcal{S}_{\theta}^{(k)}(\mathbf{x})}$. More precisely, they suggest to minimize the reverse KL divergence,

$$\begin{aligned}
 D_{\text{KL}}(\mathbf{q}_{S_{\theta}}(\mathbf{p}|\mathbf{x}) \parallel \mathbf{q}_{\mathcal{T}}(\mathbf{p}|\mathbf{x})) &= \log \Gamma \left(\sum_{k=1}^K e^{\mathcal{S}_{\theta}^{(k)}(\mathbf{x})} \right) - \sum_{k=1}^K \log \Gamma \left(e^{\mathcal{S}_{\theta}^{(k)}(\mathbf{x})} \right) \\
 &+ \sum_{k=1}^K \log \Gamma (\beta_k(\mathbf{x})) - \log \Gamma \left(\sum_{k=1}^K \beta_k(\mathbf{x}) \right) \\
 &+ \sum_{k=1}^K \left(e^{\mathcal{S}_{\theta}^{(k)}(\mathbf{x})} - \beta_k(\mathbf{x}) \right) \left(F \left(e^{\mathcal{S}_{\theta}^{(k)}(\mathbf{x})} \right) - F(\beta_k(\mathbf{x})) \right), \tag{33}
 \end{aligned}$$

where F denotes the digamma function. Throughout experiments, we also follow the practical suggestions: (1) we add one to the concentration parameters, i.e., $e^{\mathcal{S}_{\theta}^{(k)}(\mathbf{x})} \leftarrow e^{\mathcal{S}_{\theta}^{(k)}(\mathbf{x})} + 1$ and $\beta_k(\mathbf{x}) \leftarrow \beta_k(\mathbf{x}) + 1$, and (2) we minimize the loss Equation (33) divided by $\sum_{k=1}^K \beta_k(\mathbf{x})$ during optimization.

Perturbation strategies Srinivas & Fleuret (2018) show that the KD procedure on inputs perturbed by small noise implicitly encourages matching the Jacobians of a teacher and a student. One can use an isotropic Gaussian noise,

$$\tilde{\mathbf{x}} \leftarrow \mathbf{x} + \gamma \mathbf{z}, \tag{34}$$

where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \text{diag}(\mathbf{1}))$. Here, we stay consistent with the most naïve choice for the step size, that is, $\gamma = 1/255$. Throughout experiments, we adjust all perturbations to have the same size as the expected size of this Gaussian noise. Specifically, for the Output Diversified Sampling (ODS; Tashiro et al., 2020) perturbation proposed by Nam et al. (2021),

$$\tilde{\mathbf{x}} \leftarrow \mathbf{x} + \gamma \frac{\nabla_{\mathbf{x}} \mathbf{w}^{\top} \mathbf{p}_{\mathcal{T}_m}(\mathbf{x}; \tau)}{\|\nabla_{\mathbf{x}} \mathbf{w}^{\top} \mathbf{p}_{\mathcal{T}_m}(\mathbf{x}; \tau)\|_2}, \quad \text{where } \mathbf{w} \sim \mathcal{U}([-1, 1])^K, \tag{35}$$

we use the step size of $\eta = \sqrt{D}/255$, where D denotes the input dimension, i.e., $\mathbf{x} \in [0, 255]^D$. Likewise, for our perturbation strategy proposed in Section 3.2,

$$\tilde{\mathbf{x}} \leftarrow \mathbf{x} + \gamma \frac{\nabla_{\mathbf{x}}(\text{TDiv}(\mathbf{x}) - \text{SDiv}(\mathbf{x}))}{\|\nabla_{\mathbf{x}}(\text{TDiv}(\mathbf{x}) - \text{SDiv}(\mathbf{x}))\|_2}, \tag{36}$$

we use the same step size of $\eta = \sqrt{D}/255$.