
Adaptive Group Robust Ensemble Knowledge Distillation

Patrik Kenfack
ÉTS Montréal, Mila

Ulrich Aïvodji
ÉTS Montréal, Mila

Samira Ebrahimi Kahou
University of Calgary, Mila, Canada CIFAR AI Chair

Abstract

Neural networks can learn spurious correlations in the data, often leading to performance disparity for underrepresented subgroups. Studies have demonstrated that the disparity is amplified when knowledge is distilled from a complex teacher model to a relatively “simple” student model. Prior work has shown that ensemble deep learning methods can improve the performance of the worst-case subgroups; however, it is unclear if this advantage carries over when distilling knowledge from an ensemble of teachers, especially when the teacher models are debiased. This study demonstrates that traditional ensemble knowledge distillation can significantly drop the performance of the worst-case subgroups in the distilled student model even when the teacher models are debiased. To overcome this, we propose Adaptive Group Robust Ensemble Knowledge Distillation (AGRE-KD), a simple ensembling strategy to ensure that the student model receives knowledge beneficial for unknown underrepresented subgroups. Leveraging an additional biased model, our method selectively chooses teachers whose knowledge would better improve the worst-performing subgroups by upweighting the teachers with gradient directions deviating from the biased model. Our experiments on several datasets demonstrate the superiority of the proposed ensemble distillation technique and show that it can even outperform classic model ensembles based on majority voting.

1 Introduction

When trained with empirical risk minimization (ERM), neural networks are susceptible to capturing spurious correlations in the data (Tiwari and Shenoy, 2023), which are features that correlate with but not causally related to the class label (Qiu et al., 2023). In particular, the class label might spuriously correlate with patterns in the data that are easier to learn than the intended pattern. For example, in the Waterbirds dataset (Sagawa et al., 2019), which contains images of landbirds and waterbirds, most landbirds images have a land background, and waterbirds images have a water background. Instead of predicting the actual bird species, the model trained with ERM can achieve high accuracy by looking at the background. This results in a significantly higher error for underrepresented subgroups that do not exhibit the spurious correlation (e.g., landbird on water background and water bird on land background). Several works have shown that model compression methods such as pruning (Hooker et al., 2020), and knowledge distillation (Lukasik et al., 2023; Lee and Lee, 2023; Wang et al., 2023) can exacerbate the performance disparities between different subgroups. In knowledge distillation (KD), a network with a smaller capacity (student) is trained using the output of a higher capacity network (teacher) (Hinton et al., 2015). While KD can improve the student model’s average performance, the gain is not uniform across subgroups (Lukasik et al., 2023).

On the other hand, deep ensemble models have been shown to enhance generalization performance compared to individual models (Ganaie et al., 2022), and simple deep ensemble models with the same architecture, objective, and optimization settings can attenuate this shortcoming and improve the worst-case group performance (Ko et al., 2023; Kenfack et al., 2021). However, evaluating several models at test time can be computationally expensive, making them less practical for deployment on edge devices. To address this issue, ensemble knowledge distillation involves distilling the knowledge of multiple teachers to a single student model (You et al., 2017; Radwan et al., 2024), and it remains unclear whether distilling from an ensemble of teachers improves student’s worst-group performance.

This paper studies how knowledge distillation from multiple teachers impacts underrepresented subgroups. We investigate whether the subgroup performance gains of deep ensemble models apply to ensemble knowledge distillation. Focusing on logit distillation, we consider teacher models debiased by the last-layer retraining (Kirichenko et al., 2022) and investigate whether the student model can learn debiased representations from the output of retrained last-layer of the teacher. In last-layer retraining, a small held-out validation set of the group-balanced data is used to retrain the last layer of the teacher model to mitigate the spurious correlation. Our results reveal that underrepresented subgroups can be negatively impacted when distilling from multiple teachers, even when the teachers are debiased. There are other ensemble distillation approaches designed to boost the student’s performance by modeling a better aggregation of the teachers’ knowledge (Du et al., 2020; Zhang et al., 2022). In particular, the ensemble distillation method proposed by Du et al. (2020) aims to find a better compromise when teachers have conflicting predictions, and the method by Zhang et al. (2022) ensures distillation is done only using teachers with confident predictions. Our results show that these ensemble distillation methods cannot effectively fix the performance disparity of the student.

We propose an Adaptive Group Robust Ensemble Knowledge Distillation (AGRE-KD) method to encourage the student to improve the performance of unknown worst-case subgroups. Specifically, our method relies on a model that has captured the spurious correlation (i.e., a biased model) to guide the teachers’ outputs aggregation process and ensure the student model does not capture biased knowledge from the teacher. Prior work has relied on a reference classifier to target and upweight samples from unknown worst-case, using the errors of the reference classifier (Liu et al., 2021; Nam et al., 2020) or its per-sample gradient magnitude (Ahn et al., 2022). In contrast, our proposal uses the gradient direction of the biased model to select and weigh teachers’ outputs adaptively during the training.

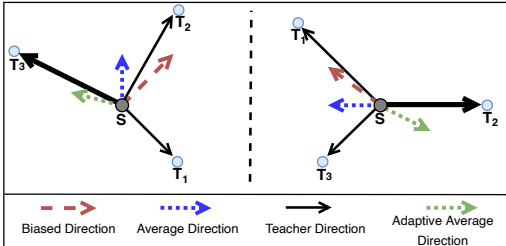


Figure 1: Illustration of our adaptive weighting method based on gradient direction. The bolder lines indicate the teacher’s higher weight in the aggregated output.

Intuitively, training the student model solely following the gradient direction that minimizes its KD loss with a biased teacher can result in local/global minima with a higher loss for the underrepresented. Our results suggest this behavior can be exacerbated in ensemble knowledge distillation since the inherent consensus from the teachers’ gradient direction can be dominated by a direction that minimizes the average error at the expense of the worst-case error. Our proposed methods compare the gradient direction of student loss with the biased model and each teacher in the ensemble and upweight teachers whose gradient direction deviates the most from the biased model since minimizing the student loss with the biased model likely results in poor worst-case group error. Figure 1 illustrates the gradient directions of the student loss with the biased model and three teachers and shows how our weighting scheme ensures the aggregated gradient direction is mainly influenced by the least biased model. As the gradient magnitude can be very noisy, we compare gradient directions by computing the dot product of their normalized vectors (i.e., the cosine similarity), which removes the influence of the gradient magnitude. The contribution of this paper can be summarized as follows:

- We demonstrate empirically that ensemble knowledge distillation amplifies performance disparities contrary to deep ensemble models. We attribute this to the reduced capacity of the student network by showing that ensemble self-distillation using models with the

same capacity reduces the performance disparity but, in some cases, achieves comparable performance with the teachers.

- We propose a novel gradient-based weighting scheme to ensure the student model minimizes teachers’ loss towards better worst-case group error. The proposed method acts in the gradient space and utilizes a model that has learned the spurious correlation (biased model) to orchestrate the distillation process.
- We perform intensive experiments on three well-known benchmarks, and the results demonstrate the superiority of the proposed method.

2 Background

We consider a multiclass classification task using a given the training $D = \{(x_i, y_i)\}_{i=1}^m$, where $x_i \in \mathcal{X}$ is input feature and $y_i \in \mathcal{Y}$ the target variable with $c = |\mathcal{Y}|$ classes, we aim to build a classifier $h(\cdot)$ that accurately predicts the class of the unlabelled test dataset. The classifier uses mapping function $f : \mathcal{X} \rightarrow \mathbb{R}^c$, that assigns scores $[\sigma_1(z_1), \dots, \sigma_c(z_c)]$, such that z is the output logits of the a given sample x and $\sigma_y(z)$ the softmax function defined as $\sigma_y(z) = \frac{\exp(z_y)}{\sum_{j \in [c]} \exp(z_j)}$, $\forall y \in [c]$.

The classifier is derived by predicting the class label that maximizes the softmax, $h(x) = \operatorname{argmax}_{j \in [c]} \sigma_j(z)$. We evaluate the classifier’s performance during training using a loss function, such as the softmax cross-entropy loss function, which measures how accurately samples are classified.

Knowledge Distillation (KD). In KD, a student network f^s aims to achieve performance close to the higher-capacity network by mimicking the teacher model f^t Hinton et al. (2015). In practice, we train the student model to mimic the teacher’s output by minimizing the Kullback-Leibler (KL) divergence between their outputs, defined as follows:

$$\mathcal{L}_{\text{KD}} = \tau^2 \cdot \text{KL}(\sigma(\frac{z^s}{\tau}), \sigma(\frac{z^t}{\tau})) \quad (1)$$

where z^s and z^t are the student and the teacher logits, respectively. τ is the temperature parameter controlling the smoothness of the probability distribution for more fine-grained information. The student loss is combined with the classification loss on the ground truth label. The overall student loss is defined by the equation 2,

$$\mathcal{L} = \alpha \cdot \mathcal{L}_{\text{KD}} + (1 - \alpha) \cdot \mathcal{L}_{\text{cls}} \quad (2)$$

where \mathcal{L}_{cls} is the classification loss (e.g., cross-entropy loss) between the student’s output and the ground truth label (y), and α is a hyperparameter controlling the classification loss and knowledge distillation loss. In other KD techniques, the student network is enforced to mimic the teacher’s internal representation instead of only outputs (Romero et al., 2014). Transferring knowledge from intermediary representation (feature-level) can provide more fine-grained information and boost the students’ performance Romero et al. (2014). Recent studies have shown that models trained with ERM still learn core features, and spurious features are only amplified in the network’s last layer (Kirichenko et al., 2022; Qiu et al., 2023; LaBonte et al., 2024). In this regard, we restrict ourselves to logit distillation and leave feature distillation for future exploration.

Like in ensemble learning, distilling knowledge from multiple teachers instead of a single one is expected to improve the student model (You et al., 2017). Ensemble learning is a widely used technique to boost the generalization performance of a model (Allen-Zhu and Li, 2020). Studies have shown that training several independent models and averaging their predictions at test time can yield a model that outperforms every individual model in the ensemble (Ganaie et al., 2022). However, evaluating multiple models for predictions at test time limits the practical use of deep ensembles due to computational overhead. Knowledge distillation can address this issue by enforcing the student network to mimic the ensemble’s output. In ensemble KD, we train the student model using the averaged softened output or the averaged knowledge distillation loss ¹ as follows:

$$\mathcal{L}_{\text{ensKD}} = \tau^2 \cdot \text{KL}(\sigma(\frac{z^s}{\tau}), \frac{1}{M} \sum_{m=1}^M \sigma(\frac{z^m}{\tau})) \quad (3)$$

¹(Du et al., 2020) showed that the averaged softened output is equal to averaged KD loss.

Where M is the number of teachers and the loss 3 can be plugged in equation 2 for training the student with KD from an ensemble of M teachers. In contrast to these methods, our method is fully unsupervised, both based in terms of ground truth class label and group information. Related works more aligned with our approach are unsupervised methods that leverage teachers’ diverse knowledge to improve students’ performance. (You et al., 2017; Du et al., 2020; Fukuda et al., 2017; Zhang et al., 2022; Kwon et al., 2020). For example, Fukuda et al. (2017) suggest that randomly selecting a teacher during mini-batch training can allow the student model to capture complementary knowledge of teachers. Other authors argue that simply averaging teachers’ softmax outputs can mislead the student model, particularly when there is competition or contradictions between teachers. In this regard, Zhang et al. (2022) proposes a sample-wise weighting for teacher loss based on the confidence of the teacher’s prediction compared to the ground label. Other methods consider label-free weighing schemes by comparing teachers in the gradient space. For instance, Du et al. (2020) use a multi-objective optimization approach in the gradient space to find teachers that agree the most in the gradient direction that minimizes their loss with the student model. Similarly, Zhou et al. (2024) sample-wise teachers selection by only averaging the majority of teachers with the same gradient directions. Our method also operates within gradient space to enhance students’ resilience to spurious correlation. Our experiments demonstrate that adhering to the majority of teachers’ opinions does not always benefit the underrepresented subgroups.

3 Related Work

Bias mitigation without group label. We consider settings where samples in the dataset D are associated with *unknown* group labels that spuriously correlate with the class label. For example, in the CelebA dataset, the class ‘hair color’ (blond, non-blond) correlates with the gender (male, female) since most images with blond hair belong to the female group. Neural networks can capture this correlation, resulting in poor performance for certain subgroups (e.g., blond males) (Sagawa et al., 2019). We aim to ensure that the model does not capture spurious correlations in the data and accurately classifies samples from all subgroups. In particular, we measure the model’s bias using the performance of *worst-case group*. Several methods have been proposed to mitigate these biases in single models (Kenfack et al., 2024b). When the group information is known, Sagawa et al. (2019) propose Group Distributionally Robust Optimization (DRO) that minimizes loss of the group experiencing the maximum loss. However, group information can be costly to collect or unavailable due to privacy restrictions (Lahoti et al., 2020; Kenfack et al., 2024a). Several methods have been proposed in this setting to improve the worst-case group performance without group labels (Kenfack et al., 2024b). For instance, Lahoti et al. (2020) proposes to use an adversary to up-weight regions where the model makes the most mistakes and demonstrates that this adversary can upweight sample from the worst-case group. Liu et al. (2021) and Nam et al. (2020) rely on a reference classifier to a reference classifier to target and upweight worst-performing. These methods use the mistakes of the reference classifier to improve group robustness by up-weighting misclassified samples (Liu et al., 2021). The reference classifier is generally trained to amplify the misclassification of samples from the unknown worst-performing group (Nam et al., 2020). While these methods do not use group labels during the training, they require access to a small validation set with group labels for model selection or hyperparameter tuning (Kenfack et al., 2024b). Kirichenko et al. (2022) proposed *Deep Feature Reweighting* (DRF) for training group robust model by training the model empirical risk minimization (ERM) on the training dataset, and then retraining the last layer of neural network with a small subset of a held-out group-balanced dataset. The proposed method achieves state-of-the-art worst-case group performance. The study demonstrated that neural networks can encode relevant and spurious features during training. Still, the impact of spurious features is upweighted in the classification layer due to high group imbalance. In a subsequent study, LaBonte et al. (2024); Qiu et al. (2023) demonstrated comparable performance by fine-tuning the last layer with fewer group information or without group labels using proxy group information from a reference classifier. This work builds on Kirichenko et al. (2022); LaBonte et al. (2024) for debiasing the teacher models by retraining their last layers using a held-out group balanced with fewer group annotations. We investigate whether distilling knowledge from an ensemble of (debiased) teachers can lead to more robust student models.

Bias in Knowledge Distillation. Several works have studied bias in knowledge distillation with a single teacher model (Lukasik et al., 2023; Lee and Lee, 2023; Lukasik et al., 2023; Tiwari et al.,

2024; Bassi et al., 2024). In particular, Lukasik et al. (2023) demonstrates that teacher errors can be amplified by the student during distillation and proposed a mitigation strategy that distills only the confident predictions of the teacher. However, their study focuses on worst-class errors and KD with a single teacher, while we study worst-subgroup errors with multiple teachers. Lee and Lee (2023) propose an adapted version of *Simple knowledge distillation* (SimKD (Chen et al., 2022)) that transplants the last layer of the teacher to student and only distills features. With the teacher trained with Group DRO, they show that transplanting the teacher’s last layer to the student only improves the worst-case group performance if the feature distillation is performed by upweighting misclassified sample from a reference classifier (Lee and Lee, 2023). However, their method uses the reference classifier by Liu et al. (2021), and the efficiency of sample weighting requires intensive tuning of the number of epochs to train the reference classifier. Similarly, Tiwari et al. (2024) uses the earlier layers of the neural network to train a reference classifier and shows that it improves the recall of worst group samples within the misclassification set, which are upweighted in the KD loss. However, their method also requires class labels to derive the misclassified samples.. Lukasik et al. (2023) study where it is best to apply the debiasing mechanism (Group DRO) and conclude that applying the robust loss to both the teacher and student model improves the average performance along with the worst-case group performance. In contrast to these prior works, we study bias in knowledge distillation with multiple teachers without group information and class labels. We investigate whether the subgroup’s performance gain observed in deep ensemble models also applies when the knowledge of the ensemble is distilled to a single model. We aim to achieve better worst-case performance across subgroups when aggregating the outputs of multiple teachers in knowledge distillation. To the best of our knowledge, this represents the first study on bias in ensemble knowledge distillation.

4 Adaptive Group Robust Ensemble Knowledge Distillation (AGRE-KD)

In this work, we considered each teacher in the ensemble to have the same architecture and trained using different random initializations. Prior work has shown ensemble models with different random initializations are diverse enough to improve the performances (Ganaie et al., 2022). This work shows that while deep ensemble models can improve the worst-case group (Ko et al., 2023), it is not necessarily the case in ensemble knowledge distillation. To address this problem, we propose AGRE-KD, an adaptive ensembling knowledge distillation strategy that ensures the student model captures robust knowledge from the teachers. Figure 4 provides an overview of our proposed method. AGRE-KD relies on a model pretrained with ERM that captured the dataset’s spurious correlation (biased). Intuitively, suppose a student model takes gradient steps toward the direction that minimizes its KD loss with the *biased* model. In that case, the resulting student model will likely capture and even amplify the reliance on the spurious correlation, i.e., the local/global minimum in that direction likely provides the worst performance for the underrepresented subgroups. Note that an additional biased model does not add extra complexity compared to related work using a reference classifier to identify the worst-performing subgroups (Nam et al., 2020; Liu et al., 2021; Kenfack et al., 2024b). Furthermore, as we will see in the experiments, our proposed method works with any biased model trained using ERM without further modifications.

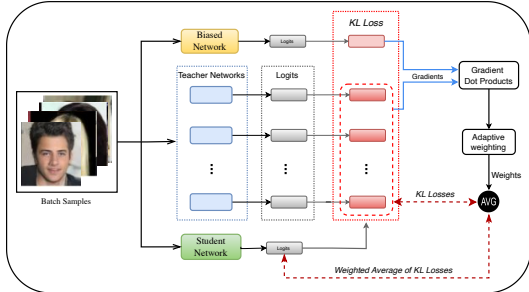


Figure 2: Overview of AGRE-KD. Detailed algorithm is provided in Supplementary A.

Gradient-based weighting scheme. For a given minibatch, we compute the student KD loss regarding the biased model b and each teacher t . Before aggregating the teachers’ outputs, we compute sample weights based on the similarity between the direction of the teacher and the biased model. Consider $\ell_i^t(\theta)$ the KD loss on sample i regarding the t -teacher, and $\ell_i^b(\theta)$ KD loss regarding the biased model, with θ the parameters of the student model. The dot product $(\langle \cdot, \cdot \rangle)$ between the *normalized*² gradients of the student KD loss with the teacher and the biased model indicates which

²Table 4 in the Supplementary shows the importance of ignoring the magnitude of the gradients.

teachers align with the biased model in gradient space. In particular, when $\langle \nabla \ell_i^t(\theta), \nabla \ell_i^b(\theta) \rangle > 0$ both models have the same gradient directions and their gradients are in the opposite direction when $\langle \nabla \ell_i^t(\theta), \nabla \ell_i^b(\theta) \rangle < 0$. In extreme cases, when the dot product of the gradients is closer to 1, both models have precisely the same directions; if the dot product is close to -1, the models are in opposite directions. Thus, to penalize teachers with a gradient direction closer to the biased model, we downweight teachers’ output for samples having dot products get closer to -1 and upweight samples as their dot products get closer to 1. We compute the sample-wise teacher’s weight as follows:

$$W_t(x_i) = 1 - \langle \nabla \ell_i^t(\theta), \nabla \ell_i^b(\theta) \rangle \quad (4)$$

Adaptive knowledge distillation. The weighting scheme in Equation 4 suggests that teacher models who behave similarly to the biased model in the gradient space will have less influence on the aggregated outputs. Therefore, we aggregate the teacher’s outputs using the sample-wise weighted average KD (wKD) loss defined as follows:

$$\mathcal{L}_{\text{wKD}} = \frac{W_t(x_i) \cdot \mathcal{L}_{\text{KD}}}{\sum_t W_t} \quad (5)$$

And we train the student model with the following final loss $\mathcal{L} = \alpha \mathcal{L}_{\text{wKD}} + (1 - \alpha) \mathcal{L}_{\text{cls}}$

In this work, we focused on unsupervised knowledge distillation, i.e., we set $\alpha = 1$. This means we train the student model only using the output of the teacher models that are eventually debiased. Supplementary A provides a detailed algorithm of the training process of AGRE-KD.

5 Experimental Results

In this section, we present the experimental setup and the empirical results. We compare our adaptive ensemble knowledge distillation to other methods and demonstrate our method’s effectiveness in improving the student model’s worst-case performance.

5.1 Setup

We evaluate the worst-case performance of the proposed method on three classification tasks: two from the vision domain (Waterbird and CelebA dataset) and one from the language domain (CivilComments)

- **Waterbird** (Sagawa et al., 2019; Liu et al., 2021) is a dataset of birds derived from Caltech-UCSD Birds (CUB) (Wah et al., 2011) by synthetically creating a spurious correlation between bird species and the background. In particular, the class label is the type of bird appearing in the image (waterbirds, landbirds), and the background landscape (water, land) spuriously correlates with the bird type. Here, the minority subgroups represent images with the background landscape not aligned with the bird type, i.e., {waterbird, land background} and {landbird, water background}.
- **CelebA** (Liu et al., 2015) dataset contains images of celebrities with 40 facial attributes. In this dataset, the attribute `hair_color` is spuriously correlated `gender`. We consider hair color {blond, non-blond} as the class label and gender {male, female} as group information.
- **CivilComments** (Koh et al., 2021) is a textual dataset collected from online comments. The task is to predict whether a comment is `toxic` or `non-toxic`. The label is spuriously correlated with comments related to some demographic subgroups such as gender (male, female), race (white, black), and sexual orientation (LGBT). We consider a binary indicator of comments related to these demographic subgroups as spurious group information.

Network Architecture and Training. Following prior work (Tiwari et al., 2024; Lee and Lee, 2023), we use the Resnet-18 (He et al., 2016) architecture for the student model and the Resnet-50 (He et al., 2016) architecture for the teacher models for the forvision tasks. Both networks are pretrained on the ImageNet-1K (Russakovsky et al., 2015) dataset. For the language task, we used the BERT (Devlin et al., 2019) model for teachers and the DistilBERT (Sanh et al., 2019) for the student

model; and language models are pretrained on Book Corpus and English Wikipedia. Following related works on KD (Du et al., 2020; Fukuda et al., 2017; Chen et al., 2022), we set the temperature hyperparameter $\tau = 4$ (Eq. 1) and show in an ablation study in Supplementary C (Figure 5) that increasing τ can exert positive effect on WGA up to certain values. We provide further details about hyperparameters in the Supplementary B.

Teacher training. We train each teacher model in the ensemble independently using standard ERM and cross-entropy loss with the ground truth class labels. Teacher models have the same architecture and hyperparameters and only differ in random seeds used for weights initialization; prior studies have shown that independent training with different random weight initializations can provide models sufficient diversity to improve the performance of the model ensemble (Allen-Zhu and Li, 2020; Ganaie et al., 2022). As aforementioned, we obtain debiased teachers using the approach introduced by Kirichenko et al. (2022) to retrain the last layer of the ERM model on the small proportion of the held-out group-balanced dataset. We perform the retraining step of the last layer using group-balanced batch sampling instead of the averaging models trained over group-balanced subsets of the data (LaBonte et al., 2024). This debiasing process is very simple and computationally inexpensive since it involves training a logistic regression model on a smaller dataset.

For the *biased* model used by our method to compute teachers’s weights, we randomly select one teacher model trained with ERM and without DFR. We provide in Supplementary C (Table 3) experiment with different choices of biased model showing the robustness of the proposed gradient-based weighting to the biased model choice. Following Kirichenko et al. (2022); LaBonte et al. (2024), we use half of the validation set of each benchmark to perform last-layer retraining with DFR; as we do not use the group and class labels during the KD training, we do not perform further hyperparameter tuning or model selection.

Baselines In addition to the standard training process using the one-hot class label for training each teacher model (Section 2), we consider other ensemble knowledge distillation methods aiming to improve the student’s performance. In particular, we consider the following baseline:

- **Deep Ensemble:** This baseline corresponds to deep ensembling using a majority voting scheme of models with the same capacity as the student model. In particular, given a set of models, the predicted class label represents the class that received the most votes for models in the ensemble.
- **KD with averaged teachers outputs (AVER)** (You et al., 2017): Here, we perform standard knowledge distillation following equation 3 by minimizing the KL loss between the student’s softmax output and the averaged softened outputs (dark knowledge) from teachers.
- **Random** (Fukuda et al., 2017): During each mini-batch training, this method randomly selects a teacher model from the ensemble to train the student model. Fukuda et al. (2017) referred to this technique as *switched-training* as the weights of the student model are updated by switching across teacher labels at the minibatch level.
- **AE-KD** (Du et al., 2020): This method is an adaptive ensembling distillation technique closest to ours. However, the method postulates that when teachers have conflicting gradient directions, a multi-objective optimization problem is solved to select the gradient direction (teachers) that satisfies most of the teachers in the ensemble.

5.2 Results

We train each model using three random seeds and report the means and standard deviations. We consider ensemble distillation with ten teacher models randomly sampled from a poll of pretrained (*biased*) teachers across random seeds. In addition to the teacher’s performance and deep ensemble approach, we report the performance of the "student" trained only using the ground truth class labels (One-hot), i.e., without knowledge distillation. We report the overall average and Worst-Group Accuracy (WGA). In the Supplementary C (Table 5), we provide the group-wise accuracy of each baseline method on the Waterbirds and CelebA datasets. We consider models trained using ERM or with last-layer retraining for debiasing. Specifically, for ensemble knowledge distillation methods, we train the student model with *biased teachers* (trained using ERM) and *debiased teachers* (last layer retrained with DFR (Kirichenko et al., 2022)). Table 1 summarizes the main results of the paper from which we draw the following observations:

Table 1: **Comparison of ensemble KD methods.** We report the average and worst-group test accuracy (WGA) on each dataset. The *Debiased* column indicates whether the model involves debiasing with DFR or whether teacher models are debiased when the *KD* column is checked (\checkmark). Bolded represents the best-performing student with the debiased teacher ensemble and underlined represents the best-performing with the biased teacher ensemble.

Models	Debiased	KD	Waterbirds		CelebA		CivilComment	
			Average	WGA	Average	WGA	Average	WGA
Teacher	\times	\times	85.7 \pm 1.80	65.6 \pm 5.75	95.4 \pm 0.07	37.5 \pm 2.27	90.0 \pm 0.25	75.9 \pm 1.15
	\checkmark	\times	94.2 \pm 0.53	90.9 \pm 1.00	93.7 \pm 2.44	90.1 \pm 1.68	85.9 \pm 0.49	77.9 \pm 0.49
Deep Ensemble	\times	\times	84.4 \pm 0.00	59.3 \pm 0.44	95.6 \pm 0.02	37.2 \pm 0.00	90.6 \pm 0.02	76.1 \pm 0.10
	\checkmark	\times	93.5 \pm 0.17	90.0 \pm 0.56	92.4 \pm 0.18	90.3 \pm 0.55	85.9 \pm 0.17	76.6 \pm 0.09
One-hot	\times	\times	83.8 \pm 0.97	54.1 \pm 2.21	95.5 \pm 0.07	32.3 \pm 2.66	90.1 \pm 0.21	75.5 \pm 0.84
	\checkmark	\times	91.4 \pm 1.31	86.7 \pm 1.85	92.3 \pm 0.39	88.9 \pm 2.90	85.7 \pm 0.47	76.5 \pm 1.09
Random	\times	\checkmark	80.0 \pm 0.37	39.6 \pm 3.63	95.2 \pm 0.03	27.2 \pm 0.00	90.8 \pm 0.21	<u>75.2</u> \pm 0.95
	\checkmark	\checkmark	89.6 \pm 0.71	77.3 \pm 2.23	92.5 \pm 0.29	85.1 \pm 0.84	91.0 \pm 0.18	74.2 \pm 1.04
AVER	\times	\checkmark	79.2 \pm 0.80	46.4 \pm 2.39	95.5 \pm 0.05	28.8 \pm 0.78	<u>90.9</u> \pm 0.03	74.7 \pm 1.16
	\checkmark	\checkmark	90.8 \pm 2.44	82.9 \pm 1.23	92.4 \pm 0.25	83.4 \pm 0.45	90.8 \pm 0.07	75.0 \pm 0.73
AE-KD	\times	\checkmark	81.7 \pm 0.80	46.9 \pm 7.57	95.3 \pm 0.06	30.5 \pm 1.88	90.8 \pm 0.54	73.1 \pm 3.05
	\checkmark	\checkmark	90.9 \pm 1.72	85.0 \pm 1.23	92.3 \pm 0.26	87.5 \pm 1.17	90.7 \pm 0.23	74.8 \pm 1.11
AGRE-KD (Ours)	\times	\checkmark	82.2 \pm 1.37	55.0 \pm 5.47	95.4 \pm 0.04	37.6 \pm 0.78	89.3 \pm 3.65	74.7 \pm 3.00
	\checkmark	\checkmark	91.3 \pm 0.49	87.9 \pm 1.23	91.7 \pm 0.20	91.9 \pm 0.71	90.2 \pm 0.49	75.9 \pm 1.75

- **The deep ensemble using models with the same architecture as the student model improves the worst-case group performance.** While for all ensemble knowledge distillation, the worst-case group performance can drop up to 10%. The drop is more severe for the random ensemble distillation but improved average accuracy, which suggests switching between teachers during the training harms worst-case groups. On the other hand, AE-KD (Du et al., 2020) performs better than other standard KD methods, showing that addressing disagreement between teachers’ gradient direction can benefit the worst group samples.
- **The fairness property of the last-layer retraining of neural network using DFR is transferable.** When we train the teacher models using ERM and then retrain the classification layer with the group-balanced set, the resulting student models achieve significantly better worst-case group accuracy. We only train the student model using the combined teachers’ outputs without any feature distillation. These results demonstrate that the re-trained classifications can provide pseudo-label distribution (dark knowledge) that reduces students’ reliance on spurious features. This shows that by mimicking the teacher’s outputs, the classifier layer of the student model also downweights the spurious features in its last layers (Kirichenko et al., 2022). However, the improved students’ WGA across ensemble distillation methods does not match the WGA of the debiased teacher models. We attribute this to the smaller capacity of the student model, which we discuss in the next experiment.
- **Our weighting scheme consistency improves the WGA compared to other ensembling methods.** On the CelebA dataset, we achieve better teachers’ average performance in terms of worst-case group accuracy as well as the deep ensemble of models with the same capacity as the student. This observation is consistent with settings where teacher models are biased or debiased. On the other hand, we can see that our method does not provide significant improvement for the worst-case group when all teacher models in the ensemble are biased; we attribute this to the fact all these teachers will likely have a similar gradient direction with biased models for many samples.

Impact of the model capacity. In this experiment, we study whether the student’s network capacity is a source of the implication of spurious feature learning in ensemble KD. In particular, we perform the same experiments as previously but using self-distillation, i.e., we use the same architecture for the student and teacher models (i.e., resnet50); we report the results on Waterbirds and Celeba datasets in Table 2 and provide results with resnet34 student in the Supplementary C (Figure 4). The ensemble KD methods with self-distillation significantly improve the worst-case test group accuracy. The gap between the teacher models and students is reduced compared to KD with a smaller capacity student

Table 2: **Results on self-distillation.** The student and the teacher models have the same network architecture (resnet50). The student’s worst-group test accuracy increases when we distill to a network with higher capacity.

Models	Debiased	KD	Waterbirds		CelebA	
			Average	WGA	Average	WGA
Deep Ensemble	✗	✗	84.4 \pm 0.00	59.3 \pm 0.44	95.6 \pm 0.02	37.7 \pm 0.00
	✓	✗	93.5 \pm 0.17	90.0 \pm 0.56	92.4 \pm 0.18	88.8 \pm 0.55
One-hot	✗	✗	85.7 \pm 1.80	65.6 \pm 5.75	95.4 \pm 0.07	37.5 \pm 2.27
	✓	✗	92.2 \pm 0.53	90.9 \pm 1.00	92.5 \pm 0.34	88.2 \pm 1.68
Random	✗	✓	83.5 \pm 0.68	64.0 \pm 2.88	95.5 \pm 0.00	35.8 \pm 1.17
	✓	✓	92.8 \pm 0.49	90.3 \pm 0.23	92.6 \pm 0.26	87.9 \pm 0.64
AVER	✗	✓	83.5 \pm 0.75	63.7 \pm 2.72	95.6 \pm 0.02	35.3 \pm 1.95
	✓	✓	92.8 \pm 0.64	90.2 \pm 0.23	92.6 \pm 0.20	88.3 \pm 1.66
AE-KD	✗	✓	84.2 \pm 1.52	61.0 \pm 4.45	95.6 \pm 0.02	36.9 \pm 0.39
	✓	✓	91.9 \pm 1.89	89.0 \pm 3.59	92.3 \pm 0.04	89.4 \pm 1.11
AGRE-KD	✗	✓	84.9 \pm 1.40	66.3 \pm 4.76	94.5 \pm 3.22	39.2 \pm 0.78
	✓	✓	91.4 \pm 1.99	91.1 \pm 2.56	91.1 \pm 0.09	91.9 \pm 1.17

model (i.e., resnet18 in Table1). The results indicate that the students’ higher capacity can help the network learn more core features and reduce the influence of spurious features in the last layer. On the other hand, AGRE-KD outperforms other baselines, showing that our adaptive weighting scheme effectively guides the student models to focus on minimizing worst-case group error during training. We further illustrate the effectiveness of our weighting scheme in the next experiment by adjusting the number of debiased teachers in the ensemble.

Effect of the number of debiased teachers in the ensemble. We study how the WGA of the student model is impacted when the teacher ensemble contains biased and debiased models with different proportions. We consider five teachers in this experiment and use the same training process as previously with different proportions of debiased teachers (i.e., $\{\frac{1}{5}, \frac{2}{5}, \dots, \frac{5}{5}\}$); we report the average and WGA in Figure 3 for the Waterbirds and CelebA datasets. The results indicate our method can adaptively identify and rely more on the knowledge of debiased teachers while reducing reliance on biased teachers. AGRE-KD can significantly improve worst-case performance when the ensemble contains a single debiased teacher. The aggregation process of other ensemble KD methods relies more on biased teachers, leading the student model to capture spurious correlation. As the proportion of debiased teachers increases, the worst-case group accuracy of all ensemble methods also increases. Additionally, AGRE-KD outperforms or matches the performance of the deep ensembling model, where the last layer of each model in the ensemble is directly retrained. These findings highlight the effectiveness of using gradients to steer model training towards specific goals, such as bias mitigation.

6 Conclusion

In this paper, we studied bias in ensemble knowledge distillation (KD) and demonstrated that, unlike deep ensemble models that reduce bias, traditional ensemble KD methods can amplify it. We proposed AGRE-KD, an adaptive gradient-based weighting method that improves group robustness in ensemble KD by guiding the student model to learn core features and boosting worst-group accuracy (WGA). Our experiments across several benchmarks demonstrated the effectiveness of our approach in distilling knowledge with reduced spurious correlations. While our results highlight AGRE-KD’s advantages over existing methods, several questions remain. First, this study focused on unsupervised KD using teachers’ logits alone, without access to group or class labels. Second, WGA improvements are less pronounced when all teachers in the ensemble are biased, suggesting an opportunity to exploit class labels in this setting to further boost WGA, e.g., by considering the teachers’ misclassifications. Finally, additional evaluations on more complex datasets, such as language tasks, are needed to validate the approach across more diverse applications.

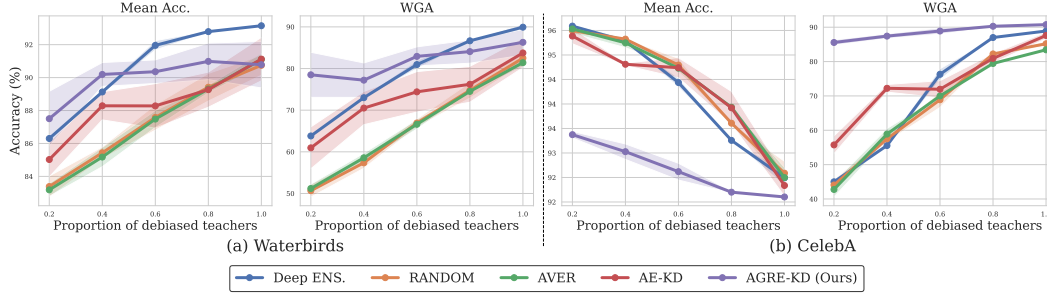


Figure 3: **Results on the proportion of debiased teachers in the ensemble.** We trained the student model using an ensemble of 5 teachers with different debiased teacher ratios within the ensemble ($\{\frac{1}{5}, \frac{2}{5}, \dots, \frac{5}{5}\}$). AGRE-KD effectively upweights and favors the least biased teachers in the ensemble, while other ensemble methods rely on them and decrease the WGA. AGRE-KD maintains significantly higher WGA, despite having only a single debiased model in the ensemble..

Acknowledgement

The authors thank the Digital Research Alliance of Canada for computing resources. SEK is supported by CIFAR and NSERC DG (2021-4086) and UA by NSERC DG (2022-04006).

References

- Ahn, S., Kim, S., and Yun, S.-Y. (2022). Mitigating dataset bias by using per-sample gradient. *arXiv preprint arXiv:2205.15704*.
- Allen-Zhu, Z. and Li, Y. (2020). Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*.
- Bassi, P. R., Cavalli, A., and Decherchi, S. (2024). Explanation is all you need in distillation: Mitigating bias and shortcut learning. *arXiv preprint arXiv:2407.09788*.
- Chen, D., Mei, J.-P., Zhang, H., Wang, C., Feng, Y., and Chen, C. (2022). Knowledge distillation with the reused teacher classifier. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11933–11942.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*.
- Du, S., You, S., Li, X., Wu, J., Wang, F., Qian, C., and Zhang, C. (2020). Agree to disagree: Adaptive ensemble knowledge distillation in gradient space. *advances in neural information processing systems*, 33:12345–12355.
- Falcon, W. and team, T. P. L. (2019). Pytorch lightning. Available at: <https://github.com/Lightning-AI/lightning>.
- Fukuda, T., Suzuki, M., Kurata, G., Thomas, S., Cui, J., and Ramabhadran, B. (2017). Efficient knowledge distillation from an ensemble of teachers. In *Interspeech*, pages 3697–3701.
- Ganaie, M. A., Hu, M., Malik, A., Tanveer, M., and Suganthan, P. (2022). Ensemble deep learning: A review. *Engineering Applications of Artificial Intelligence*, 115:105151.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Hinton, G., Vinyals, O., and Dean, J. (2015). Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.

- Hooker, S., Moorosi, N., Clark, G., Bengio, S., and Denton, E. (2020). Characterising bias in compressed models. *arXiv preprint arXiv:2010.03058*.
- Izmailov, P., Kirichenko, P., Gruver, N., and Wilson, A. G. (2022). On feature learning in the presence of spurious correlations. *Advances in Neural Information Processing Systems*, 35:38516–38532.
- Kenfack, P. J., Kahou, S. E., and Aïvodji, U. (2024a). Fairness under demographic scarce regime. *Transactions on Machine Learning Research*.
- Kenfack, P. J., Kahou, S. E., and Aïvodji, U. (2024b). A survey on fairness without demographics. *Transactions on Machine Learning Research*.
- Kenfack, P. J., Khan, A. M., Kazmi, S. A., Hussain, R., Oracevic, A., and Khattak, A. M. (2021). Impact of model ensemble on the fairness of classifiers in machine learning. In *2021 International conference on applied artificial intelligence (ICAPAI)*, pages 1–6. IEEE.
- Kirichenko, P., Izmailov, P., and Wilson, A. G. (2022). Last layer re-training is sufficient for robustness to spurious correlations. *arXiv preprint arXiv:2204.02937*.
- Ko, W.-Y., D’souza, D., Nguyen, K., Balestrierio, R., and Hooker, S. (2023). Fair-ensemble: When fairness naturally emerges from deep ensembling. *arXiv preprint arXiv:2303.00586*.
- Koh, P. W., Sagawa, S., Marklund, H., Xie, S. M., Zhang, M., Balsubramani, A., Hu, W., Yasunaga, M., Phillips, R. L., Gao, I., et al. (2021). Wilds: A benchmark of in-the-wild distribution shifts. In *International conference on machine learning*, pages 5637–5664. PMLR.
- Kwon, K., Na, H., Lee, H., and Kim, N. S. (2020). Adaptive knowledge distillation based on entropy. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7409–7413. IEEE.
- LaBonte, T. (2023). Milkshake: Quick and extendable experimentation with classification models. <http://github.com/tmlabonte/milkshake>.
- LaBonte, T., Muthukumar, V., and Kumar, A. (2024). Towards last-layer retraining for group robustness with fewer annotations. *Advances in Neural Information Processing Systems*, 36.
- Lahoti, P., Beutel, A., Chen, J., Lee, K., Prost, F., Thain, N., Wang, X., and Chi, E. (2020). Fairness without demographics through adversarially reweighted learning. *Advances in neural information processing systems*, 33:728–740.
- Lee, J. and Lee, J. (2023). Debaised distillation by transplanting the last layer. *arXiv preprint arXiv:2302.11187*.
- Liu, E. Z., Haghighi, B., Chen, A. S., Raghunathan, A., Koh, P. W., Sagawa, S., Liang, P., and Finn, C. (2021). Just train twice: Improving group robustness without training group information. In *International Conference on Machine Learning*, pages 6781–6792. PMLR.
- Liu, Z., Luo, P., Wang, X., and Tang, X. (2015). Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738.
- Loshchilov, I., Hutter, F., et al. (2019). Fixing weight decay regularization in adam. *International Conference on Learning Representations (ICLR)*.
- Lukasik, M., Bhojanapalli, S., Menon, A. K., and Kumar, S. (2023). Teacher’s pet: understanding and mitigating biases in distillation. *Transactions on Machine Learning Research*.
- Mohammadshahi, A. and Ioannou, Y. (2024). What is left after distillation? how knowledge transfer impacts fairness and bias. *arXiv preprint arXiv:2410.08407*.
- Nam, J., Cha, H., Ahn, S., Lee, J., and Shin, J. (2020). Learning from failure: De-biasing classifier from biased classifier. *Advances in Neural Information Processing Systems*, 33:20673–20684.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. (2017). Automatic differentiation in pytorch.

- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. (2019). Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.
- Qiu, S., Potapczynski, A., Izmailov, P., and Wilson, A. G. (2023). Simple and fast group robustness by automatic feature reweighting. In *International Conference on Machine Learning*, pages 28448–28467. PMLR.
- Radwan, A., Zaafarani, L., Abudawood, J., AlZahrani, F., and Fourati, F. (2024). Addressing bias through ensemble learning and regularized fine-tuning. *arXiv preprint arXiv:2402.00910*.
- Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. (2014). Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. (2015). Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252.
- Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P. (2019). Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*.
- Sanh, V., Debut, L., Chaumond, J., and Wolf, T. (2019). Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *ArXiv*, abs/1910.01108.
- Tiwari, R. and Shenoy, P. (2023). Overcoming simplicity bias in deep networks using a feature sieve. In *International Conference on Machine Learning*, pages 34330–34343. PMLR.
- Tiwari, R., Sivasubramanian, D., Mekala, A., Ramakrishnan, G., and Shenoy, P. (2024). Using early readouts to mediate featural bias in distillation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2638–2647.
- Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. (2011). The caltech-ucsd birds-200-2011 dataset. *California Institute of Technology*.
- Wang, S., Narasimhan, H., Zhou, Y., Hooker, S., Lukasik, M., and Menon, A. K. (2023). Robust distillation for worst-class performance: on the interplay between teacher and student objectives. In *Uncertainty in Artificial Intelligence*, pages 2237–2247. PMLR.
- You, S., Xu, C., Xu, C., and Tao, D. (2017). Learning from multiple teacher networks. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1285–1294.
- Zhang, H., Chen, D., and Wang, C. (2022). Confidence-aware multi-teacher knowledge distillation. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4498–4502. IEEE.
- Zhou, W., Ding, Z., Zhang, X., Shi, H., Wang, J., and Yin, D. (2024). Govern: Gradient orientation vote ensemble for multi-teacher reinforced distillation. *arXiv preprint arXiv:2405.03764*.

A Algorithm

Algorithm 1 Adaptive Group Robust Ensemble Knowledge Distillation (AGRE-KD)

- 1: **Input:** Ensemble of pretrained teachers $T = \{T_1, T_2, \dots, T_M\}$, biased model T_b , student model S with parameters θ , dataset \mathcal{D} , distillation coefficient α , temperature parameter τ
- 2: **for** each minibatch $\{(x_i, y_i)\}_{i=1}^B$ from \mathcal{D} **do**
- 3: **for** each teacher $T_t \in T$ **do**
- 4: Compute knowledge distillation loss $\ell_i^t(\theta)$ for sample i between S and T_t ▷ (Eq. 1).
- 5: Compute biased model distillation loss $\ell_i^b(\theta)$ for sample i between S and T_b ▷ (Eq. 1).
- 6: Compute gradient alignment $G_i^t = \langle \nabla \ell_i^t(\theta), \nabla \ell_i^b(\theta) \rangle$ ▷ Dot product of normalized gradient vectors. Table 4 shows the importance of ignoring the magnitude of the gradients.
- 7: Compute the adaptive sample weight for teacher t on sample i : $W_t(x_i) = 1 - G_i^t$;
- 8: **end for**
- 9: Compute weighted knowledge distillation loss:

$$\mathcal{L}_{\text{wKD}} = \frac{\sum_t W_t(x_i) \cdot \ell_i^t(\theta)}{\sum_t W_t(x_i)}$$

- 10: Compute classification loss (if labeled data available): \mathcal{L}_{cls}
 - 11: Compute final loss:

$$\mathcal{L} = \alpha \mathcal{L}_{\text{wKD}} + (1 - \alpha) \mathcal{L}_{\text{cls}}$$
 - 12: Update student model parameters θ using gradient descent on \mathcal{L}
 - 13: **end for**
-

B Hyperparameters

We train all models using standard hyperparameters from previous work (LaBonte et al., 2024; Kirichenko et al., 2022; Izmailov et al., 2022) and keep their value fixed across experiments. For the vision tasks, we used an initial learning rate of 1×10^{-3} with a cosine learning rate scheduler; we used a batch size of 32 and 100 for the Waterbirds and the CelebA datasets, respectively. For the CiviComments dataset, we use an initial learning rate of 1×10^{-5} with a linear learning rate scheduler, a batch size 16, and train for ten epochs. We keep all hyperparameters fixed to train the teacher and student models. For the optimizer, we used AdamW (Loshchilov et al., 2019) and SGD for the language and vision datasets, respectively, with a weight decay of 1×10^{-4} . Our implementation uses PyTorch (Paszke et al., 2017, 2019), Torch Lightning (Falcon and team, 2019), and Milkshake (LaBonte, 2023).

C Supplemental experiments

Table 3: **Sensitivity of AGRE-KD to the biased model architecture.** In the main paper, we used a biased model with the same architecture as the teacher models (i.e., resnet50). We experiment with different network backbones for the biased models in AGRE-KD. The results below do not show significant differences across biased model choices, demonstrating the robustness of the proposed method to the choice of biased model. These results suggest that the gradient direction of any biased pretrained model can provide sufficient guidance for debiased distillation.

Biased model	Waterbirds		CelebA	
	Average	WGA	Average	WGA
Resnet50	90.6 \pm 0.49	86.7 \pm 2.82	91.8 \pm 0.20	90.5 \pm 0.24
Resnet34	90.0 \pm 1.63	85.2 \pm 3.48	91.8 \pm 0.14	90.5 \pm 0.24
Resnet18	90.2 \pm 1.10	86.8 \pm 2.88	91.5 \pm 0.09	89.8 \pm 0.16

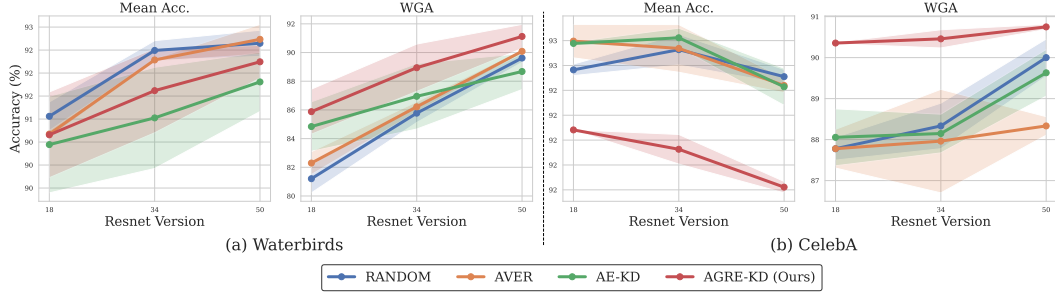


Figure 4: **Results on student model capacity.** We perform experiments on WaterBirds and CelebA using resnet50 teachers and varying the capacity of the student network (resnet18, resnet34, and resnet50). We plot the average and the worst-case accuracy of different ensemble distillation methods across three random seeds. The WGA tends to increase as the student model has more capacity for learning the core features. Most baselines match or outperform the teachers when the student model has the same capacity as the teachers (i.e., resnet50), and our method remains superior in terms of WGA. These results suggest that the reduced capacity of the student model is a source of the disparity observed.

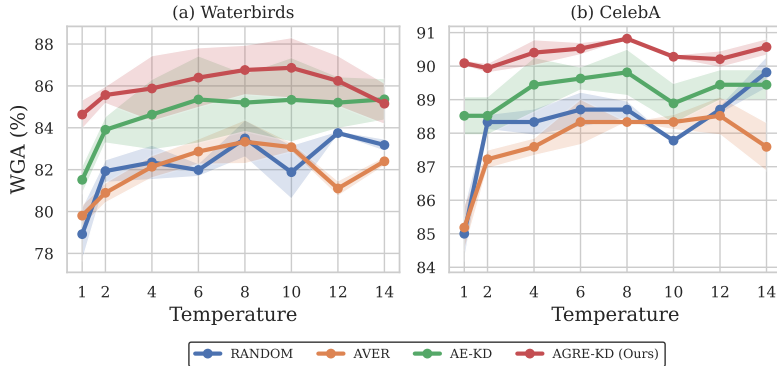


Figure 5: **Effect of the temperature parameter on the WGA.** We study the effect of the temperature parameter (Eq. 3) in WaterBirds and CelebA datasets using five resnet50 teachers and a resnet18 student. We perform knowledge distillation with different temperature parameters $\tau \in \{1, 2, 4, \dots, 10\}$ and report the average test WGA across three independent random seeds. The WGA can improve with increased temperature and steady or decrease after certain higher values of temperature ($\tau \geq 10$ and $\tau \geq 8$ for the Waterbirds and CelebA datasets, respectively). These results align with recent work by Mohammadshahi and Ioannou (2024) studying the benefit of increased temperature for fairness in knowledge distillation.

Table 4: **Training AGRE-KD using gradient direction with (w/) and without (w/o) gradient normalization.** On the Waterbirds and CelebA datasets, we study how using the gradient magnitude in the weighting scheme impact the results. We trained our AGRE-KD method without normalizing the gradient vectors in the dots product, i.e., accounting for the gradient magnitude of the losses. The results below show that accounting for the gradient magnitude in the weighting scheme reduces the WGA performance. This shows the importance of using normalized dot products to compare directional changes in gradients, making the computed weights independent of the gradient scales themselves, which are generally very noisy.

Model	Waterbirds		CelebA	
	Average	WGA	Average	WGA
AGRE-KD w/ grad norm	90.6 \pm 0.49	86.8 \pm 1.86	91.7 \pm 0.20	90.9 \pm 0.71
AGRE-KD w/o grad norm	91.0 \pm 1.32	81.5 \pm 4.07	93.1 \pm 0.21	87.9 \pm 0.52

Table 5: **Group-wise accuracy comparison.** We report the group-wise accuracy of different ensemble KD methods on the Waterbirds and CelebA datasets. As in previous experiments, we average the performances across three different random seeds and considered ensembles of five teachers.

Datasets	Sub Groups	#Samples	ERM	Teacher	Random	AVER	AEKD	AGRE-KD
Waterbirds	(landbirds,land)	3498	99.3 \pm 0.13	95.5 \pm 0.97	96.7 \pm 0.57	96.8 \pm 0.62	96.1 \pm 1.29	94.3 \pm 1.46
	(landbirds,water)	184	74.0 \pm 2.54	94.3 \pm 0.96	86.2 \pm 2.20	87.4 \pm 0.91	86.4 \pm 3.88	87.5 \pm 3.14
	(waterbirds,land)	56	54.1 \pm 1.80	92.1 \pm 1.39	82.3 \pm 1.65	82.1 \pm 1.06	84.6 \pm 3.37	86.7 \pm 2.82
	(waterbirds,water)	1057	93.4 \pm 0.72	90.9 \pm 0.81	91.5 \pm 1.08	90.9 \pm 0.67	91.4 \pm 1.97	90.6 \pm 1.85
CelebA	(nonblond,female)	71629	96.0 \pm 0.39	91.1 \pm 0.37	91.8 \pm 0.52	91.8 \pm 0.52	91.3 \pm 0.36	90.4 \pm 0.77
	(nonblond,male)	66874	99.9 \pm 0.06	92.9 \pm 0.20	94.3 \pm 0.37	94.2 \pm 0.20	93.7 \pm 0.40	92.4 \pm 0.53
	(blond,female)	22880	85.2 \pm 1.21	94.4 \pm 0.30	93.8 \pm 0.35	93.9 \pm 1.13	94.6 \pm 0.34	95.0 \pm 0.83
	(blond,male)	1387	32.3 \pm 2.17	90.1 \pm 0.86	88.3 \pm 0.78	87.5 \pm 0.52	89.4 \pm 1.63	91.6 \pm 0.78

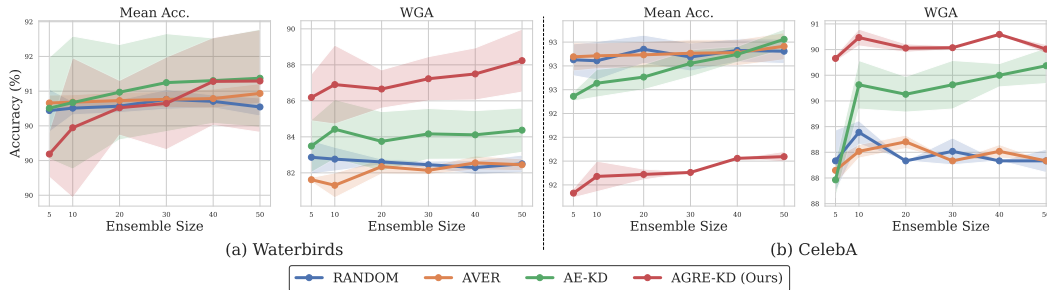


Figure 6: **Effect of teacher size.** We train the resnet18 student model with different ensemble sizes of resnet50 teachers (5, 10, 20, ..., 50). Increasing ensemble size exerts a positive effect on both the average and the worst-group performance. However, the Random KD method tends to get worse as we increase the number of teachers.