General Greedy De-bias Learning

Xinzhe Han, Student Member, IEEE, Shuhui Wang, Member, IEEE, Chi Su, Qingming Huang, Fellow, IEEE, and Qi Tian, Fellow, IEEE

Abstract—Neural networks often make predictions relying on the spurious correlations from the datasets rather than the intrinsic properties of the task of interest, facing with sharp degradation on out-of-distribution (OOD) test data. Existing de-bias learning frameworks try to capture specific dataset bias by annotations but they fail to handle complicated OOD scenarios. Others implicitly identify the dataset bias by special design low capability biased models or losses, but they degrade when the training and testing data are from the same distribution. In this paper, we propose a General Greedy De-bias learning framework (GGD), which greedily trains the biased models and base model. The base model is encouraged to focus on examples that are hard to solve with biased models, thus remaining robust against spurious correlations in the test stage. GGD largely improves models' OOD generalization ability on various tasks, but sometimes over-estimates the bias level and degrades on the in-distribution test. We further re-analyze the ensemble process of GGD and introduce the Curriculum Regularization inspired by curriculum learning, which achieves a good trade-off between in-distribution (ID) and out-of-distribution performance. Extensive experiments on image classification, adversarial question answering, and visual question answering demonstrate the effectiveness of our method. GGD can learn a more robust base model under the settings of both task-specific biased models with prior knowledge and self-ensemble biased model without prior knowledge. Codes are available at https://github.com/GeraldHan/GGD.

Index Terms—Dataset Biases, Robust Learning, Greedy Strategy, Curriculum Learning

1 INTRODUCTION

EEP learning have been used in a wide range of tasks that involves vision and/or language [1]. Most of the current approaches are data-driven and heavily rely on the assumption that the training and testing data are drawn from the same distribution. They are usually susceptible to poor generalization on out-of-distribution or biased settings [2]. This limitation partially arises because supervised training only identifies the correlations between given examples and their labels [3], which may reflect the dataset-specific bias rather than intrinsic properties of the task of interests [4], [5]. In general, under the supervised objective function fitting paradigm, if the bias is sufficient to make the model achieve high accuracy, there is less motivation for models to further learn those true instrinsic factors of the task. For example, QA models trained on SQuAD [6] tend to select the text near question-words as answers regardless of the context [7], [8], and VQA models usually leverage superficial correlations between questions and answers without considering the vision information [9], [10]. When it comes with the more common situation that the distribution of test data deviates from that of training data, models exploiting the biases in

training data are prone to show poor generalization and hardly provide proper evidence for their predictions.

Being aware of this problem, researchers re-examine many popular datasets, resulting in the discovery of a wide variety of biases on different tasks, such as language bias in VQA [11], color bias in Biased-MNIST [12], gender/background bias in image classification [13], [14], and the ubiquitous long-tailed distribution [15], [16]. Built on these findings, explicit de-bias methods [12], [17], [18], [19], [20], [21] assume that bias variables are explicitly annotated, then the out-of-distribution performance can be directly improved by preventing the model from using the known biases or bias-related data augmentation [22], [23]. Although these methods achieve remarkable improvement on typical diagnosing datasets, they can only mitigate one specific bias, which is inconsistent with the real world datasets with compositional biases [24]. For example, in VQA, biases may stem from unbalanced answer distribution, spurious language correlations, and object contexts. Even when all bias variables are identified, the explicit de-bias methods still cannot well handle multiple types of biases. Some recent works, *i.e.*, the implicit methods [25], [26], [27], [28], try to discover the compositional biases without explicit taskrelated prior knowledge. They are somehow overcomplicated and usually perform worse than explicit methods under welldefined circumstances with known biases.

In fact, the dataset biases can be reduced in a more straightforward manner. As shown in Fig. 1, features learned from biases are thought to be "spurious" because they can only generalize to the majority groups of samples in the dataset. Although the model may incur high training error on the minority groups where the spurious correlation does not hold, the overall loss will still be trapped in a local minimum due to the low average training error dominated by the majority groups. Compared with the core features

Corresponding author: Shuhui Wang.

[•] X. Han and Q. Huang are with the School of Computer Science and Technology, University of Chinese Academy of Sciences, Beijing 101408, China, and with the Key Laboratory of Intelligent Information Processing, Institute of Computing Technology, Chinese Academy of Sciences, Beijing 100190, China. Q. Huang is also with Peng Cheng Laboratory, Shenzhen 518066, China.

E-mail: xinzhe.han@vipl.ict.ac.cn, qmhuang@ucas.ac.cn.

S. Wang is with the Key Laboratory of Intelligent Information Processing, Institute of Computing Technology, Chinese Academy of Sciences, Beijing 100190, China, and with Peng Cheng Laboratory, Shenzhen 518066, China. E-mail: wangshuhui@ict.ac.cn.

C. Su is with SmartMore, Beijing, 100085. Email: chi.su@smartmore.com

Q. Tian is with Huawei Cloud & AI, Shenzhen 518129, China. E-mail: tian.gi1@huawei.com.



Fig. 1. Examples of dataset biases in different tasks. Models tend to capture spurious correlations between the inputs and the labels instead of the task of interest. From top to bottom, we illustrate the color bias in Biased-MNIST, background bias in image classification, and compositional language bias in VQA.

(*e.g.*, the semantics of objects), it is relatively easier to identify the biases brought by distractive information (*e.g.*, the backgrounds). If it is possible to know ahead of time that which subsets of instances are irrelevant to spurious features, we can encourage the model to focus on these samples and then reduce the unexpected correlations. In our preliminary work [29], we propose a de-bias framework Greedy Gradient Ensemble (GGE) to mitigate the language biases in VQA, achieving great improvement on the biased dataset VQA-CP [11]. GGE greedily learns a series of biased models and then ensembles the biased models and the base model like gradient descent in the functional space. The gradient of biased model naturally indicates the difficulty of a sample with certain spurious correlation.

However, without explicit labels for bias variables, even when the prior knowledge for dataset bias is given, disentangling the biased part from the whole feature still remains ill-posed [30]. Degradation on in-distribution test is a common problems in exiting de-bias method [18], [22], [25] including GGE. Greedily emphasizing the bias reduction would also lead to the overreach of the learning objectives. Ever worse, this bias overestimation brings harm to the model generalizablity under more general cases. Few works have considered this issue for de-bias learning [31], [32], relying on extra effort in constructing a biased model separately for ensemble or distillation, but they appear to be less flexible in dealing with complex real applications. In this paper, we first re-analyse the cause of bias overestimation in GGE. We empirical find that if a large amount of data can be correctly predicted via the biased model with high confidence, they will be excluded in the training of base model. As a result, the base model may be under-fitted to

some labels due to inadequate training data. Decomposing the negative gradient of cross-entropy loss, we further find that the cross-entropy between the base prediction and biased prediction measures the difficulty of samples in GGE.

Base on this finding, we transform the negative gradient supervision to a flexible regularization and formulate a more general framework, *i.e.*, General Greedy De-bias (GGD), to tackle the bias over-estimation problem more appropriately. Inspired by the curriculum learning [33], we treat the regularization term as a difficulty metric for the curriculum selection function. In this way, all data can participate in base model training in the early stage and gradually focus on hard examples along with the training procedure. This treatment endows our model with more flexibility, and demonstrates robustness on both out-of-distribution test data and general datasets like ImageNet [34] and CIFAR [35].

In the experiments, we apply GGD to a wider range of uni-modal and multi-modal tasks, including visual classification, linguistic question answering, and visual question answering. Quantitative and qualitative evaluations on all the tasks show that our framework is feasible to general dataset bias on different tasks and gains improvement on both in-distribution (ID) and out-of-distribution (OOD) performance without extra annotations in training and extra computational cost in inference.

The main contributions of this paper are summarized as:

- We present a de-bias framework, General Greedy Debias Learning, which encourages unbiased based model learning by the robust ensemble of biased models. GGD is more generally applicable compared to task-related explicit de-bias learning methods while more flexible and effective compared to implicit de-bias methods.
- We propose Curriculum Regularization for GGD, which results in a new training scheme GGD_{cr} that can better alleviate the "bias over-estimation" phenomenon. Compared with previous methods [36], [37], GGD_{cr} comes to a better trade-off between in-distribution and out-of-distribution performance without either extra unbiased data in training or model ensemble in inference.
- Experiments on image classification, question answering, and visual question answering demonstrate the effectiveness of GGD on different types of biases.

This paper provides a more general debias learning framework compared to our preliminary study [29]. First, the previous work [29] only aims at the VQA task, in which we pay attention to the bias analysis on VQA-CP [11] and new evaluation metric for models' visual grounding ability, while this paper considers the general de-bias learning problem and extend our framework to various datasets and applications. Second, we provide discussions for the "bias over-estimation" phenomenon in previous GGE [29]. We propose a flexible GGD_{cr} optimization scheme that effectively improves the in-distribution performance on different tasks. Third, we provide more in-depth analysis for the greedy de-bias strategy, and the differences between GGD and previous GGE are demonstrated from both theories and experiments. Fourth, we provide more experiments on GGD with known biases and unknown biases, and comparison with the latest de-bias methods with both explicit and implicit bias modelling are also provided. Finally, we apply

GGD on three additional tasks, *i.e.*, image classification, adversarial question answering, and visual question answering, which are the representative tasks from CV, NLP, and Vision-Language, respectively, demonstrating the circumstances under one single bias, unknown bias, long-tailed bias and multiple biases. In addition to VQA-CP and VQA v2 in [29], we add experiments on more datasets, *i.e.*, Biased MNIST [12], SQuAD [6], Adversarial SQuAD [38], GQA-OOD [39], CIFAR-10 and CIFAR-100 [40]. Experiments on various tasks and datasets fully demonstrate the general applicability of GGD in debias learning.

2 RELATED WORK

2.1 De-biasing from Data Sources

When collecting real-world datasets, biases in the data are inevitable. Torralba and Efros [2] show how biases affect some commonly used datasets. It draws consideration on the generalization performance and classification capability of the trained deep models. Recent dataset construction protocols have tried to avoid certain kinds of biases. For example, on both CoQA [41] and QuAC [42] for QA task, annotators are prevented from using words that occur in the context passage. For VQA, Zhang *et al.* [43] collect complementary abstract scenes with opposite answers for all binary questions. Similarly, VQA v2 [9] is introduced to weaken the language priors in the VQA v1 dataset [44] by adding similar images with different answers for each question.

However, constructing large-scale datasets is costly. It is crucial to develop models that are robust to biases [45]. Towards this goal, new diagnosing datasets are established by amplifying some specific biases. For instance, Agrawal *et al.* [11] constructed a diagnosing VQA dataset under Changing Prior (VQA-CP), with different answer distributions between the train and test splits. Adversarial SQuAD [38] is built by adding distracting sentences to the passages in SQuAD [6]. He *et al.* collect NICO [14] dataset that consists of images with different backgrounds and gestures. All these new datasets can be used to test the models' generalization ability on out-of-distribution scenarios.

2.2 Explicitly De-biasing with Known Bias

To train a de-biased model, some works utilize an intentionally biased model to de-bias another model. For VQA, Ramakrishnan et al. [46] introduce an adversarial regularization to remove the discriminative features related to the answer categories from the questions. RUBi [18] and PoE [47] re-weight samples based on the question-only predictions. Kim et al. [20] propose a regularization term based on mutual information between the feature embedding and the bias, to remove the known bias for image classification. Similarly, Clark et al. [17] construct bias-only models for VQA, reading comprehension, and natural language inference (NLI), then reduce them with bias production and entropy maximization. Xiong et al. [30] further conduct uncertainty calibration on the bias-only models for a better de-biasing performance. It can detect sample outliers and feature noises simultaneously. Bahng et al. [12] find that Hilbert-Schmidt Independence Criterion (HSIC) [48] can encourage a set of features to be

statistically independent. They capture local texture bias in image classification and static bias in the video action recognition task using small-capacity models and then train a de-biased representation that is independent of biased representations based on HSIC.

Teney et al. [23] generate counterfactual samples with specific prior knowledge for different tasks. The vector difference between pairs of counterfactual examples serves to supervise the gradient orientation of the network. Liang et al. [49] propose A-INLP that dynamically finds bias-sensitive tokens and mitigates social bias in text generation. Tartaglione et al. [21] propose a new regularization named EnD, which aims to disentangle the features having the same "bias label". Sagawa et al. [19] avoid bias over-fitting by defining prior data sub-groups and controlling their generalization. HEX [50] pushes the model to learn representations from which the texture representation is not predictable with the reverse gradient method. Gat et al. [51] introduce a regularization by maximizing functional entropies (MFE), which forces the model to use multiple information sources in multimodal tasks. Zhu et al. [52] explicitly extract target and bias features from the latent space. Then they learn to discover and remove their correlation with the mutual information estimation. Hong et al. [53] leverage the knowledge of bias labels and propose Bias-Contrastive and Bias-Balanced losses based on the contrastive learning.

The above methods only focus on one specific bias but cannot work well on compositional biases. GGD can sequentially mitigate multiple bias variables as long as they can be characterized with prior knowledge, which is much more flexible than explicit de-biasing methods.

2.3 Implicitly De-biasing without Known Bias

In real-world scenario, bias presented in the dataset is often hard to characterize and disentangle. To address this issue, there have been several recent works to resolve dataset bias without explicit supervision on the biases. For linear models, to alleviate the co-linearity among variables, Shen *et al.* [54] propose to learn a set of sample weights that can make the design matrix nearly orthogonal. Kuang *et al.* [55] further propose a re-weighting strategy so that the weighted distribution of treatment and confounder could satisfy the independent condition.

For deep models, most implicit methods assume that easy-to-learn biases can be captured by models with limited capacity and model parameters [56], using a small subset of training instances in a few epochs [57], and a classifier attached to intermediate layers [58]. Apart from limited capacity biased models, Huang et al. [27] iteratively discard the dominant features activated on training data and force the network to activate the remaining features correlated with labels. Nam et al. [25] amplify the biases using generalized cross-entropy (GCE) loss and train a de-biased classifier with resampling based on the biased classifier. Still based on GCE, BiaSwap [28] further generates bias-swapped images from bias-contrary images as bias-guided data augmentation. Zhang et al. [59] introduce a non-linear feature decorrelation approach based on Random Fourier Features, which can approximate the Hilbert-Schmidt norm in Euclidean space. Spectral Decoupling [26] decouples the learning dynamics

between features. It aims to overcome the issue of gradient starvation, which indicates the tendency to only rely on statistically dominant features. Moreover, the setting of implicit de-biasing is similar to the Domain Generalization (DG) [60], [61] but has different challenges. In DG, the model is encouraged to generalize to a new domain that is not accessible during training while the de-bias has a small amount of training data that is bias-conflicted. Meanwhile, since there is no clear "domain discrepancy" in the biased sets, most existing DG methods do not work on the dataset bias problem.

Implicit methods are much more flexible. However, compared with explicit de-bias methods, totally ignoring prior knowledge limits their capability upper-bound for some tasks. If multiple types of biases are characterized, they cannot fully leverage all the valuable information. In contrast, GGD makes use of task-specific knowledge so that it can mitigate compositional biases. For tasks without prior knowledge of the biases, it can also learn a more robust model with a self-ensemble biased model like implicit de-biasing methods, gaining more flexibility in real world applications.

3 PROPOSED METHOD

3.1 Preliminaries

In this section, we first introduce the notations used in the rest of this paper. $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ denotes the training set, where \mathcal{X} is the feature space of observations, and \mathcal{Y} is the label space. Assume $\mathcal{B} = \{B_1, B_2, \ldots, B_M\}$ to be a set of task-specific bias features that can be extracted in priority, such as texture features in Biased-MNIST and the language shortcut in VQA. Correspondingly, $h_m(B_m; \phi_m) : B_m \to \mathcal{Y}$ is a biased model that makes prediction with certain biased feature B_m , where ϕ_m is the parameter set of $h_m(.)$ that maps B_m to the label space \mathcal{Y} . Similarly, $f(X; \theta) : \mathcal{X} \to \mathcal{Y}$ denotes the base model, *i.e.*, our target model for inference. For supervised learning, the training objective is to minimize the distance between the predictions and the labels Y as

$$\min_{\theta} \mathcal{L}\left(f(X;\theta),Y\right),\tag{1}$$

where the loss function can be various types of supervision loss, such as cross-entropy (CE) loss for single-label classification, binary cross-entropy (BCE) loss for multilabel classification, triplet loss for retrieval, *etc.* Similar to previous works [12], [21], [25], [26], considering that the classification (and its variants) is the most common task that seriously suffers from the dataset bias problem, we also take classification tasks as a demonstration in this paper.

3.2 Greedy Gradient Ensemble

Given Eq. 1, f(.) is chosen to be an over-parametrized DNN, so the model is easy to over-fit the biases in the datasets and suffers from poor generalization ability. We take advantage of the easy-to-overfit property of deep models, and joinly fit the ensemble of bias models $\sum_{m=1}^{M} h_m(B_m; \phi_m)$ and base model $f(X; \theta)$ to label Y

$$\min_{\phi,\theta} \mathcal{L}\left(f(X;\theta) + \sum_{m=1}^{M} h_m(B_m;\phi_m), Y\right).$$
 (2)

Algorithm 1: GGD_{qs} **Input:** Observations *X*, Labels *Y*, Biased feature Observations $\mathcal{B} = \{B_m\}_{m=1}^M$, Base function $f(.|\theta): X \to \mathbb{R}^{|Y|}$, Bias functions $\{h_m(.|\phi_m): B_m \to \mathbb{R}^{|Y|}\}_{m=1}^M$ Initialize: $\mathcal{H}_0 = 0$; for Batch $t = 1 \dots T$ do for $m = 1 \dots M$ do $L_m(\phi_m) \leftarrow$ $\mathcal{L}'(h_m(B_m;\phi_m), -\nabla \mathcal{L}(H_{m-1}, Y))$ Update $\phi_m \leftarrow \phi_m - \alpha \nabla_{\phi_m} L_m(\phi_m)$ end $L_{M+1}(\theta) \leftarrow \mathcal{L}'(f(X;\theta), -\nabla \mathcal{L}(H_M, Y))$ Update $\theta \leftarrow \theta - \alpha \nabla_{\theta} L_{M+1}(\theta)$ end return $f(X;\theta)$

Ideally, we hope the spurious correlations are *only* over-fitted by the bias models, thus the base model f(.) can be learned with a relatively unbiased data distribution. To achieve this goal, GGE adopts a greedy strategy that encourages biased models to have a higher priority to fit the dataset. In practice, f(.) can be ResNet for image classification, UpDn [62] for VQA, *etc.*, while h(.) can be low capability model for the texture bias, question-answer classifier for the question shortcut bias, *etc.*.

Viewing from a general ensemble model in the functional space [63], suppose we have $\mathcal{H}_m = \sum_{m'=1}^m h_{m'}(B_{m'})$ and we wish to find $h_{m+1}(B_{m+1})$ added to \mathcal{H}_m so that the loss $\mathcal{L}(\sigma(\mathcal{H}_m + h_{m+1}(B_{m+1})), Y)$ decreases. Theoretically, the desired direction of h_{m+1} should be the negative derivative of \mathcal{L} at \mathcal{H}_m , *i.e.*,

$$-\nabla \mathcal{L}(\mathcal{H}_{m,j}) := \frac{\partial \mathcal{L}(\mathcal{H}_m, Y)}{\partial \mathcal{H}_{m,j}}, j \in 1, 2, ..., C.$$
(3)

where $\mathcal{H}_{m,j}$ denotes the prediction for the *j*-th class among the overall *C* classes. For a classification task, we only care about the probability for class $j: \sigma(f_j(x)) \in (0, 1)$. Therefore, we treat the negative gradients as pseudo labels for classification and optimize the new model $h_{m+1}(B_{m+1})$ with

$$\mathcal{L}\left(h_{m+1}(B_{m+1};\phi_{m+1}),-\nabla\mathcal{L}(\mathcal{H}_m)\right).$$
(4)

After integrating all biased models, the expected base model f is optimized with

$$\mathcal{L}(f(X;\theta), -\nabla \mathcal{L}(\mathcal{H}_M)).$$
 (5)

In the test stage, we only use the base model for prediction. In order to make the above paradigm adaptive to mini-Batch Gradient Decent (MBGD), we implement an iterative optimization scheme [29] as shown in Algorithm 1. Note that our framework learns the base model and biased models jointly, which is different from existing work [32], [36] where the biased model is learned via another independent process or additional annotations.

3.3 General Greedy De-bias Learning

As shown in [29], GGD_{gs} (GGE) often over-estimates the biases of datasets. It achieves remarkable improvement on



Fig. 2. Comparison between GGD_{gs} and GGD_{cr} . GGD_{gs} (GGE) uses the gradient from the biased model as the pseudo label while GGD_{cr} enlarges the prediction discrepancy between the base model and the biased model with curriculum learning. GGD_{gs} is a special case of GGD when $\lambda_t = 1$ under CE loss.

Algorithm 2: GGD_{cr} **Input:** Observations *X*, Labels *Y*, Biased feature Observations $\mathcal{B} = \{B_m\}_{m=1}^M$, Base function $f(.|\theta): X \to \mathbb{R}^{|Y|}$, Bias functions $\{h_m(.|\phi_m): B_m \to \mathbb{R}^{|Y|}\}_{m=1}^M$ Initialize: $\mathcal{H}_0 = 0$; for Batch $t = 1 \dots T$ do $\lambda_t \leftarrow \sin(\frac{\pi t}{2T})$ for $m = 1 \dots M$ do $L_m(\phi_m) \leftarrow$ $\begin{array}{l} \mathcal{L}'\left(h_m(B_m;\phi_m),-\nabla\mathcal{L}(H_{m-1},Y)\right)\\ \text{Update }\phi_m\leftarrow\phi_m-\alpha\nabla_{\phi_m}L_m(\phi_m) \end{array}$ end $\hat{\sigma}(\mathcal{H}_M) \leftarrow Y \odot \sigma(\mathcal{H}_M)) \ L_{M+1}(\theta) \leftarrow$ $\mathcal{L}\left(f(X;\theta),Y\right) - \lambda_t CE(f(X),\hat{\sigma}(\mathcal{H}_M))$ Update $\theta \leftarrow \theta - \alpha \nabla_{\theta} L_{M+1}(\theta)$ end return $f(X;\theta)$

out-of-distribution data but may significantly degrades under the in-distribution setting. To overcome this critical issue, we first re-analyse the biased model in GGE under CE loss

$$\mathcal{L}_{CE}(Z,Y) = -\sum_{j=1}^{C} y_j \log(\sigma_j),$$
(6)

with

$$\sigma_j = \frac{e^{z_j}}{\sum_{k=1}^C e^{z_k}},\tag{7}$$

where $Z = \{z_j\}_{j=1}^C$ is the predicted logits, and $y_j \in \{0, 1\}$ is the ground-truth label for the *j*-th class. σ_j indicates the confidence of the biased model on *j*-th class. The negative gradient of the loss function is

$$-\nabla \mathcal{L}(z_j) = y_j - \sigma_j. \tag{8}$$

To make the range of pseudo labels consistent with the classification label space [0,1], $-\nabla L(z_j)$ is clipped to

$$-\nabla \hat{\mathcal{L}}(z_j) = \begin{cases} y_j - \sigma_j & y_j > 0\\ 0 & y_j = 0 \end{cases}.$$
(9)

The negative gradients access whether a sample can be solved based on the spurious correlation captured by certain biased model.

Now, casting aside the viewpoint of gradient descent in functional space, we can also decompose the CE loss with $-\nabla \hat{\mathcal{L}}$ as pseudo label

$$\mathcal{L}_{CE}(f(X), -\nabla \hat{\mathcal{L}}) = -\sum_{j=1}^{C} (y_j - \hat{\sigma}_j) \log(p_j)$$

$$= -\sum_{j=1}^{C} y_j \log(p_j) + \sum_{j=1}^{C} \hat{\sigma}_j \log(p_j)$$

$$= \mathcal{L}_{CE}(f(X), Y) - \mathcal{L}_{CE}(f(X), \hat{\sigma}),$$
(10)

where the reference prediction $\hat{\sigma} = Y \odot \sigma$ and \odot is the element-wise product that equals to the clipping in Eq. 9.

Based on Eq. 10, the gradient ensemble actually aims to provide predictions that agree with the ground-truth but disagree with the biased models. $-\mathcal{L}_{CE}(f(X), \hat{\sigma})$ controls the degree of spurious relation to be reduced. To this end, we can treat $-\mathcal{L}_{CE}(f(X), \hat{\sigma})$ as a regularization:

$$L(\theta) = \mathcal{L}\left(f(X;\theta), Y\right) - \lambda_t \mathcal{L}_{CE}(f(X), \hat{\sigma}(\mathcal{H}_M)), \quad (11)$$

where λ_t denotes the weight of the regularization term. This more general framework is noted as General Greedy De-bias (GGD), where we only keep greedy strategy but get free from the negative gradient supervision. GGE is a special case of GGD when $\lambda_t = 1$. We will denote GGE as GGD_{gs} (Gradient Supervision) in the following paper.

Furthermore, inspired by Curriculum Learning [33], $-\mathcal{L}_{CE}(f(X), \hat{\sigma})$ can be regarded as a *soft* difficulty measurement for curriculum sample selection function. In practice, we formulate a Curriculum Regularization training scheme (GGD_{cr}), which gradually increases λ_t along with the training process. In this way, samples with spurious correlations to the labels can participate in the early stage of training. In the consequent training stage, the model will focus on the hard samples that cannot be solved by biased models, resulting in more stable prediction on out-of-distribution data. The overall optimization procedure GGD_{cr} is shown in Algorithm 2. Comparison between GGD_{gs} and general GGD_{cr} is shown in Fig. 2.

3.4 Discussions

3.4.1 Intuitive Explanation of GGD

Section 3.2 has presented theoretical evidence for GGD_{gs} from the aspect of model learning in functional space. More intuitively, GGD_{gs} can also be regarded as a re-sampling strategy [64]. For a sample that is easy to fit by biased models, $-\nabla \hat{\mathcal{L}}(z_i)$ (*i.e.*, the pseudo label produced by the base model) will become relatively small. This makes $f(X;\theta)$ pay more attention to samples that are hard to fit by previous ensemble biased classifiers. As a result, the base model is not likely to learn biased features. This hard example mining process is experimentally demonstrated in Section 4.1.3 and Fig. 9.

However, according to Eq. 9, samples that demonstrate high spurious correlations (*i.e.*, $-\nabla \hat{\mathcal{L}}(z_i) = 0$) will be discarded. If large groups of data are absent because of zero supervision, the representation learning of the base model with the gradient supervision may be under-fitted. Moreover, when the label distribution is skewed (*e.g.*, distribution bias in VQA-CP [11]), the base model may over-estimate the bias in labels. This results in "inverse" training bias and significant degradation on in-distribution test data. Experiments on longtailed classification also revealed similar findings [65], [66], which indicate that re-sampling a part of the data encourages a more balanced classifier but harms the representation learning stage, while learning with unbalanced data results in a biased classifier but still provides a good representation.

To alleviate the "bias over-estimation", GGD_{cr} provides a good relaxation of GGD_{gs} by replacing the gradient supervision with a "softer" Curriculum Regularization. By adjusting λ_t , all data can participate in the base model learning in the early stage, thus the bias over-estimation can be well alleviated. We will further experimentally demonstrate these findings in Section 4.1.3 and Section 4.4.

3.4.2 Probabilistic Justification

Following the assumptions in [17], for a given sample x, let x^b be the biased features and x^{-b} be the features except the biases. x^b and x^{-b} are conditionally independent given the label y. We have

$$\log p(y|x^{-b}) = \log p(y|x) - \log p(y|x^{b}) + C, \quad (12)$$

where *C* is a constant term related to the given datasets. The detailed derivation is provided in the Appendix A. It is hard to distinguish the core features for the task of interest (x^{-b}) but it is easier to identify the dominant biases (x^b) based on the prior knowledge. Eq. 12 indicates that maximizing the likelihood $\log p(y|x^{-b})$ equals to maximizing $\log p(y|x)$ while minimizing $\log p(y|x^b)$.

Assume the optimal biased model $h(x^b; \phi^*)$ has

$$\phi^* = \arg\min_{\phi} \mathbb{E}_{\langle X, Y \rangle} \mathcal{L}(h(x^b; \phi), y).$$
(13)

Taking $q_{\phi^*}(y|x^b)$ as the distribution of optimal biased prediction $h(x^b; \phi^*)$, GGD alternatively minimizes $\log p(y|x^b)$ by enlarging the divergence between p(y|x) and the biased reference $q_{\phi^*}(y|x^b)$. Maximizing Eq. 12 is approximated as

$$\arg\max_{\theta} \left(\log p_{\theta}(y|x) + D(p_{\theta}(y|x)) || q_{\phi^*}(y|x^b)) \right)$$
(14)

where θ is the parameter of the base model that produce distribution p(y|x) and D(.||.) is the divergence between

6

two distributions. In practice, we get diverse predictions by maximizing the cross-entropy between p(y|x) and $q(y|x^b)$. Similar implementation also appears in [67]. Eq.12 provides a new justification of GGD from probabilistic formulation, which aims to maximum the log-likelihood of $\log p(y|x^{-b})$. Moreover, the precision of $q(y|x^b)$ is crucial. If the biased model captures the true correspondence too much, maximizing the divergence will harm the base model.

3.4.3 On the Trade-off between ID and OOD Performance

The key idea of greedy ensemble is similar as the Boosting strategy [68], [69]. Boosting is to combine multiple weak classifiers with high bias but low variance to produce a strong classifier with low bias and low variance. Each base learner has to be weak enough, otherwise, the first few classifiers will easily over-fit to the training data [70]. Different from boosting that ensembles all weak learners, we make use of this over-fitting phenomenon but only use the last base model for prediction. This strategy removes specific spurious correlations with the biased models but also encounters the bias-overestimation with a single base model. To solve this problem, GGD introduce Curriculum Regularization, which trains the model with all ID data in the early training stage and then gradually focus on the hard samples.

The trade-off between ID and OOD performance has already attracted much attention in the study of OOD generalization. Most of these methods assume that the OOD data is available during training [36], [67], [71], [72] or the model can be pre-trained on balanced data with few biases [37], [73]. Therefore, they can adaptively adjust the model with the given OOD data. However, for de-bias learning, the absence of OOD data makes the bias estimation more ill-posed and challenging. The works in [31], [32] share similar idea with our GGD_{cr} in that they aim to make full use of the biased ID data to pursue a good trade-off between ID and OOD performance. However, both [31] and [32] have to train a biased and a de-biased model separately and then combine the two to achieve more robust predictions. In comparison, GGD learns the two models under the unified framework as in Algorithm 2. It does not require extra training cost of an original model, and can well adapt to any choice of base model, thus it gains more flexibility in real applications.

3.5 General Applicability of GGD

This section provides the detailed instantiation of GGD on specific tasks. In the following part, let h(.) denote the biased model and $\hat{y} \in \mathcal{Y}$ denote the biased predictions, where the super-script represents the bias type.

3.5.1 GGD with Single Explicit Bias

In order to compare with existing explicit de-bias methods that focus on one single type of bias, we first test GGD on the texture bias in Biased-MNIST [12].

The dataset $\mathcal{D} = \{x_i, y_i, b_i\}_{i=1}^N$ consists of a synthetic image x_i , the annotated digit label y_i , and the background color b_i . We aim to predict the digit number \hat{y}_i with the input image x_i

$$\hat{y}_i = f(x_i),\tag{15}$$

TABLE 1

Comparison on Biased-MNIST. ρ_{train} and ρ_{test} denote the level of texture bias during training and testing, respectively. 1k uses SimpleNet-1k as the biased model, bg adopts the backgrounds as biased feature, and se stands for the self-ensemble version. 'Original' is the original MNIST without texture bias.

0		0.990		0.995			0.999			Original
$\rho_{\rm train} / \rho_{\rm test}$	0	0.1	0.990	0	0.1	0.995	0	0.1	0.999	Oligiliai
Baseline	$77.80{\scriptstyle \pm 1.30}$	$80.11{\scriptstyle \pm 0.81}$	99.76±0.06	$53.78{\scriptstyle\pm2.08}$	58.03±2.73	$99.82{\scriptstyle \pm 0.08}$	$10.44{\scriptstyle\pm2.36}$	$17.39{\scriptstyle\pm3.71}$	$99.80{\scriptstyle \pm 0.09}$	$98.78{\scriptstyle \pm 0.15}$
ReBias ^{1k} [12]	$85.21{\scriptstyle\pm0.59}$	$87.80{\scriptstyle \pm 0.62}$	$99.77{\scriptstyle \pm 0.06}$	$73.60{\scriptstyle\pm1.18}$	$75.95{\scriptstyle\pm1.58}$	$99.70{\scriptstyle \pm 0.22}$	$32.84{\scriptstyle\pm4.02}$	$37.52{\scriptstyle\pm2.37}$	$99.87{\scriptstyle \pm 0.01}$	$99.05{\scriptstyle\pm0.08}$
RUBi ^{1k} [18]	87.36±3.59	$91.30{\scriptstyle\pm2.18}$	$99.15{\scriptstyle \pm 0.29}$	$77.84{\scriptstyle \pm 6.53}$	$82.22{\scriptstyle\pm4.14}$	$99.21{\scriptstyle \pm 0.67}$	$30.25{\scriptstyle \pm 9.64}$	$37.20{\scriptstyle\pm7.90}$	$92.15{\scriptstyle \pm 7.71}$	$98.90{\scriptstyle \pm 0.04}$
GGD_{gs}^{1k}	$92.79{\scriptstyle \pm 0.76}$	$94.22{\scriptstyle\pm0.78}$	$98.69{\scriptstyle \pm 0.56}$	$91.27{\scriptstyle\pm0.25}$	$91.80{\scriptstyle \pm 0.59}$	$98.16{\scriptstyle \pm 1.00}$	$67.57{\scriptstyle\pm4.31}$	$70.77{\scriptstyle \pm 3.99}$	$86.84{\scriptstyle\pm 6.35}$	$98.64{\scriptstyle \pm 0.04}$
GGD_{cr}^{1k}	$91.78{\scriptstyle\pm1.18}$	$92.30{\scriptstyle\pm1.17}$	$99.64{\scriptstyle \pm 0.30}$	$83.90{\scriptstyle\pm2.30}$	$84.91{\scriptstyle\pm2.68}$	$99.28{\scriptstyle \pm 0.15}$	$68.36 \pm \textbf{1.89}$	$70.70{\scriptstyle\pm2.02}$	$99.25{\scriptstyle \pm 0.35}$	$99.14{\scriptstyle \pm 0.05}$
ReBias ^{bg} [12]	$84.95{\scriptstyle\pm1.63}$	$86.88{\scriptstyle \pm 1.96}$	99.66±0.25	$74.27{\scriptstyle\pm3.50}$	$76.20{\scriptstyle \pm 1.78}$	$99.74{\scriptstyle\pm0.12}$	$27.74{\scriptstyle\pm8.07}$	$34.20{\scriptstyle\pm 6.67}$	$99.87{\scriptstyle\pm0.01}$	$99.87{\scriptstyle\pm0.17}$
RUBi ^{bg} [18]	$88.65{\scriptstyle \pm 0.47}$	$89.67{\scriptstyle \pm 0.64}$	$99.59{\scriptstyle \pm 0.33}$	$78.19{\scriptstyle \pm 5.06}$	$80.50{\scriptstyle \pm 4.18}$	$98.79{\scriptstyle \pm 1.07}$	$21.07{\scriptstyle\pm6.78}$	$27.59{\scriptstyle\pm 6.39}$	$90.16{\scriptstyle \pm 4.35}$	$98.70{\scriptstyle \pm 0.06}$
GGD^{bg}_{gs}	$93.78{\scriptstyle \pm 1.34}$	$94.46{\scriptstyle \pm 1.09}$	$99.01{\scriptstyle \pm 0.42}$	$90.34{\scriptstyle \pm 0.95}$	$91.26{\scriptstyle \pm 1.06}$	$99.38{\scriptstyle \pm 0.39}$	61.81±4.29	66.00±4.77	$91.25{\scriptstyle\pm1.98}$	$98.77{\scriptstyle\pm0.07}$
GGD_{cr}^{bg}	$90.64{\scriptstyle \pm 0.84}$	$91.95{\scriptstyle\pm0.91}$	$99.82{\scriptstyle \pm 0.05}$	$86.20{\scriptstyle\pm1.41}$	$87.02{\scriptstyle\pm1.04}$	$99.68{\scriptstyle \pm 0.10}$	62.96±5.74	$67.62{\scriptstyle \pm 4.51}$	$99.41{\scriptstyle \pm 0.41}$	$99.07{\scriptstyle\pm0.13}$
ReBias ^{se} [12]	$83.77{\scriptstyle\pm0.81}$	$85.76{\scriptstyle\pm0.71}$	$99.77{\scriptstyle\pm0.08}$	$75.06{\scriptstyle \pm 3.35}$	$77.25{\scriptstyle \pm 3.61}$	$99.84{\scriptstyle\pm0.08}$	31.82±3.49	$38.41{\scriptstyle\pm2.61}$	$99.87{\scriptstyle\pm0.02}$	$99.03{\scriptstyle \pm 0.06}$
RUBi ^{se} [18]	$27.37{\scriptstyle\pm8.04}$	$33.47{\scriptstyle\pm 6.20}$	$89.14{\scriptstyle\pm8.02}$	$16.23{\scriptstyle \pm 6.95}$	$22.96{\scriptstyle\pm 5.84}$	$95.77{\scriptstyle\pm 5.18}$	$10.21{\scriptstyle \pm 5.28}$	16.67±3.79	$83.67{\scriptstyle\pm11.51}$	-
GGD^{se}_{gs}	$79.35{\scriptstyle \pm 1.53}$	$80.78{\scriptstyle\pm2.11}$	$94.65{\scriptstyle\pm 5.41}$	69.70±3.22	$72.49{\scriptstyle\pm3.13}$	$90.61{\scriptstyle \pm 1.48}$	$38.72{\scriptstyle\pm4.00}$	$42.74{\scriptstyle\pm3.25}$	$76.24{\scriptstyle \pm 3.74}$	$93.88{\scriptstyle\pm8.87}$
GGD^{se}_{cr}	$83.28{\scriptstyle\pm0.65}$	$85.53{\scriptstyle\pm1.32}$	$99.27{\scriptstyle\pm0.27}$	$72.91{\scriptstyle \pm 2.49}$	$76.19{\scriptstyle \pm 1.61}$	$99.34{\scriptstyle \pm 0.27}$	$43.78{\scriptstyle \pm 2.82}$	$48.92{\scriptstyle \pm 1.64}$	$99.46{\scriptstyle \pm 0.21}$	$98.94{\scriptstyle\pm0.05}$

where the base model f(.) is a neural network trained with CE loss.

The bias for Biased-MNIST comes from the spurious correlation between the digits and the background colors. In practice, we define two different kinds of bias models. In the first case, the biased prediction B_t^i of an image sample x_i is extracted with a low capacity model

$$\hat{y}_{i}^{t} = h_{1k}(x_{i}).$$
 (16)

 $h_{1k}(.)$ is the SimpleNet-1k [12] with kernel size 1×1 . It will predict the target class of an image only through the local texture cues due to small receptive fields.

In the second case, we provide the explicit background b_i for bias extraction

$$\hat{y}_i^t = h_{bg}(b_i). \tag{17}$$

 $h_{bg}(.)$ is a common neural network similar to the base model but the input is only a background image without digits. Therefore, the biased model will purely make predictions according to the texture bias. The experimental analysis is provided in Section 4.1.

3.5.2 GGD with Self-Ensemble

For tasks like Adversarial QA [38], the task-specific biases are hard to distinguish. For de-bias learning at the lack of prior knowledge, we design a more flexible version of GGD with Self-Ensemble, named GGD^{se} . The biased predictions B_{se} is captured with

$$\hat{y}_i^{se} = h_{se}\left(x_i\right),\tag{18}$$

where $h_{se}(.)$ is another neural network that has the same architecture and optimization scheme as the baseline model. Since the baseline model usually tends to over-fit the dataset biases, $h_{se}(.)$ can implicitly capture the biases without task-specific prior knowledge.

In the experiments, we will demonstrate the hardexample-mining mechanism of GGD^{se} on Adversarial SQuAD [38] in Section 4.2 and further verify its generalization ability on all the other three tasks.

3.5.3 GGD with Multiple Biases

To verify whether GGD can handle multiple types of biases, we conduct experiment on the Language bias in VQA. As analysed in [29], the language bias is mainly composed of two aspects, *i.e.*, distribution bias and shortcut bias.

We consider the formulation of VQA task as a classification problem. Given a dataset $\mathcal{D} = \{v_i, q_i, a_i\}_{i=1}^N$ consisting of an image $v_i \in \mathcal{V}$, a question $q_i \in \mathcal{Q}$ and a labeled answer $a_i \in \mathcal{A}$, we need to optimize a mapping $f_{VQ} : V \times Q \to \mathbb{R}^C$ which produces a distribution over the *C* answer candidates. The function is as follows

$$\tilde{a}_i = f_{\theta}(v_i, q_i) = c \left(m \left(e_v(v_i), e_q(q_i) \right) \right),$$
 (19)

where $e_v : \mathcal{V} \to \mathbb{R}^{n_v \times d_v}$ is an image encoder, $e_q : \mathcal{Q} \to \mathbb{R}^{n_q \times d_q}$ is a question encoder, m(.) stands for the multimodal fusion module, and c(.) is the multi-layer perception classifier. The output vector $\tilde{a} \in \mathbb{R}^C$ indicates the probability distribution on all the answer candidates.

The distribution bias is the statistical answer distribution under certain question types

$$\hat{y}_i^d = p(a_i|t_i),\tag{20}$$

where t_i denotes the type of question q_i , such as "what color", "is this", *etc.*, in VQA v2 [44].

The shortcut bias is the semantic correlation between specific QA pairs, which can be modeled as a question-only branch similar to [18]

$$\hat{y}_i^q = c_q \left(e_q(q_i) \right), \tag{21}$$

where $c_q: Q \to \mathbb{R}^C$.

To verify whether GGD can handle compositional biases, we design different versions of GGD which ensemble distribution bias, shortcut bias and both biases. The experimental results are shown in Section 4.3.



Fig. 3. Per-class Accuracies on Biased MNIST. All methods are trained with $\rho_{\text{train}} = 0.999$ and tested on $\rho_{\text{test}} = 0.1$. The upper row is the confusion matrix between the predictions and the ground-truth labels; the lower row shows the confusion matrix between the predicted labels and the background color labels.

4 EXPERIMENTS

In this section, we present experiments for GGD on both ID and OOD settings. With respect to different types of biases, experiments on image classification [12], QA [38] and VQA [44] are shown afterwards, corresponding to CV, NLP and Vision-Language tasks. Note that GGD is a general debias framework, and it will be an interesting issue in further study to apply our methods on other tasks that also suffer from dataset biases.

4.1 Image Classification

4.1.1 Dataset

Biased MNIST. To better analyse the the properties of GGD, we first verify our model on Biased MNIST, where we can have full control over the amount of bias during training and evaluation. Biased MNIST [12] is modified from MNIST [74] which introduces the color bias that highly correlates with the label *Y* during training.

On Biased-MNIST, 10 different colors are selected for each digit $y \in \{0, ..., 9\}$. For each image of digits y, we assign a pre-defined color with probability ρ and any other color with probability $1 - \rho$. $\rho \in [0, 1]$ controls the level of spurious correlation in train and test set. $\rho = 0.99$ means 99% images in the dataset are assigned with a background of the corresponding color. $\rho = 0.1$ is the unbiased condition with a uniform sampled background color.

4.1.2 Experimental Setups

For image classification on biased MNIST, we use ResNet-18 as our baseline model. Model 1k in TABLE 1 denotes SimpleNet-1k with kernel size 1×1 proposed in [12]. Biased model *bg* uses the ResNet-18 with background color images as the input. All experiments use the same CE loss and baseline model. "Original" stands for MNIST dataset without biases. Considering the randomness in Biased-MNIST data generation, we report the mean and variance of 4 repeated experiments under different random seeds.

4.1.3 Experimental Results

GGD Overcomes Bias. As shown in TABLE 1, GGD largely improves the OOD Accuracy on Biased MNIST. Under extremely biased training data ($\rho_{\text{train}} = 0.999$), the best performed method $\mathrm{GGD}_{cr}^{\bar{1}k}$ achieves 68% accuracy on the unbiased test data ($\rho_{\text{test}} = 0$), which is over 6 times compared with the baseline model. In Fig. 3, we provide per-class accuracy matrix for more detailed analysis. The diagnostic heat-map corresponds to unbiased and biased correlation respectively. We can observe that the vanilla ResNet-18 mainly captures the spurious correlation but confuses on the ground-truth digits. ReBias [12] can better capture the core correlations compared with the baseline but will be still fooled by the texture bias. GGD hardly relies on the background color as shown in Fig. 3 (c) and (d), which demonstrate that our method can help a model to overcome certain kinds of bias via specially designed biased model.

Performance on Different Bias Level. As shown in TA-BLE 1, GGD achieves prominent performance gain on out-ofdistribution tests across all bias levels while remaining stable on in-distribution data. Comparing with other methods that use the ensemble strategy, GGD surpasses RUBi [18] and ReBias [12] by a large margin under the same base model and biased models, especially when the training and testing data are extremely different ($\rho_{\text{train}} = 0.999$). When training and testing under the unbiased situation, both GGD_{cr} and GGD_{gs} are stable if the bias type is known ahead. RUBi^{se} fails on unbiased training set, which even continuously decreases under the original MNIST ('-' in Table 1).

Ablations on Biased Models. Besides biased model with small receptive field, we test another version of the biased model with a ground-truth background image as the biased

TABLE 2 Ablations for λ_t . "Anneal" indicates the Curriculum Regularization that changes λ_t from 0 to 1 along with training process.

λ.	Train 0.999						
Λ_t	0	0.1	0.999				
0 (Baseline)	8.95	16.24	99.87				
$1 (GGD_{gs})$	68.31	71.34	91.42				
0.95	58.79	63.41	99.96				
Anneal (GGD $_{cr}$)	67.01	70.17	99.58				

TABLE 3

Experimental Results on Adversarial QA. We provide the F1 score on Adversarial SQuAD AddSent split and SQuAD v1 dev split. "Original" is trained with SQuAD train split, and "Extra" is trained with extra Adversarial SQuAD AddSentOne split.

Method -	Orig	ginal	Extra			
	AddSent	Dev	AddSent	Dev		
Baseline	$46.61{\scriptstyle \pm 0.30}$	$87.61{\scriptstyle \pm 0.16}$	$50.11{\scriptstyle \pm 0.35}$	$87.72{\scriptstyle \pm 0.21}$		
GGD^{se}_{gs}	$48.05{\scriptstyle \pm 0.11}$	$87.98{\scriptstyle \pm 0.38}$	$52.44{\scriptstyle \pm 0.49}$	$87.01{\scriptstyle\pm0.61}$		
GGD^{se}_{cr}	$48.42{\scriptstyle\pm0.20}$	$87.89{\scriptstyle \pm 0.20}$	$53.94{\scriptstyle\pm0.09}$	$88.38{\scriptstyle \pm 0.27}$		

features, shown as *bg* in TABLE 1. GGD works well with different biased models comparing with other methods.

For implicit de-biasing, a biased model with the same structure as the baseline ResNet-18 is trained in the Self-Ensemble version GGD^{se} . As shown in TABLE 1, GGD^{se}_{cr} achieves the best performance when the bias information is not available. RUBi^{se} corrupts under the self-ensemble setting. This demonstrates that GGD^{se} can also implicitly remove biases even without the task-specific biased models, which is much more flexible compared with existing explicit de-bias methods.

 GGD_{gs} vs. GGD_{cr} . As shown in TABLE 1, although achieving high Accuracy on the out-of-distribution test data, GGD_{gs} is not as robust as GGD_{cr} with the increase of texture bias. Especially on the in-distribution test data, the accuracy is significantly lower than the baseline ResNet-18. Moreover, GGD_{gs}^{se} is also very unstable in the later training stage, resulting in large variance according to different training data. On the other hand, GGD_{cr}^{se} achieves comparable indistribution performance against the baseline even on the original MNIST dataset.

For better analysis of the GGD_{gs} and GGD_{cr}, we design another ablation study on λ_t in Eq. 11 under $\rho_{\text{train}} = 0.999$. As shown in TABLE 2, by slightly relaxing the regularization (changing λ_t from 1.0 to 0.95), the in-distribution accuracy will be increased to the level of vanilla ResNet-18. This verifies our assumption that such degradation mainly comes from the *completely* absence of samples with the spurious correlation (Section 3.3). Starting from this insight, we define $\lambda_t = \sin(\frac{\pi t}{2T})$ in GGD_{cr}, where t is the current training epoch and T is the number of total epochs. With this Curriculum Regularization, GGD_{cr} achieves a good tradeoff between in-distribution and out-of-distribution tests, remaining comparable in-distribution test accuracy against the baseline and OOD test against GGD_{gs}.

4.2 Adversarial Question Answering

4.2.1 Dataset

For the NLP tasks, we choose the adversarial question answering (AdQA) to demonstrate the effectiveness of GGD. We evaluate on the Adversarial SQuAD [38] dataset, which was built by adding distractive sentences to the passages in SQuAD [6]. These sentences are designed to be very similar to the corresponding questions but with a few key semantic changes to ensure that they do not indicate the correct answer. Models that only focus on the similarity between question and context will tend to be misled by the new sentence. A sample from Adversarial SQuAD is shown in Fig. 4.

• Question: What is the name of the quarterback who was 38 in Super Bowl XXXIII?
Context: Peyton Manning became the first quarterback ever to
lead two different teams to multiple Super Bowls. He is also the
oldest quarterback ever to play in a Super Bowl at age 39. The past
record was held by John Elway, who led the Broncos to victory in
Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice
President of Football Operations and General Manager. Jeff Dean is
the name of the quarterback who was 37 in Champ Bowl XXXIV.
• GT Answer: John Elway
Misleading Answer: Jeff Dean

Fig. 4. An example from Adversarial SQuAD. Blue sentence is the expected evidence, while red sentence is the distracting sentence to fool the models.

4.2.2 Experimental Setups

We use BiDAF [81] as the base model. It introduces a multi-stage hierarchical process that represents the context at different levels of granularity and uses a bidirectional attention flow mechanism to obtain a query-aware context representation without early summarization. Since the word-level and semantic-level similarity is hard to distinguish by modelling, we only test GGD^{se} in the following experiments.

The models are trained and validated on the original SQuAD train and val set, and test on the AddSent split of Adversarial SQuAD. In order to further verify the hard example mining mechanism behind GGD, we also design another "Extra" setting in which we add the AddSentOne split of Adversarial SQuAD to the training set as additional hard samples. The performances are measured with F1 score, which are the weighted average of the precision and recall rate at the character level.

4.2.3 Experimental Results

The F1 scores for the OOD test on Adversarial SQuAD and ID test on SQuAD v1 are shown in TABLE 3. The reported results are from four repeated experiments with different random seeds. We find that both GGD_{gs} and GGD_{cr} , trained with the SQuAD train set, only improve the performance by ~ 2%. However, after adding a few hard examples from AddSentOne, both methods achieve much more improvement on Adversarial SQuAD compared with the baseline. GGD_{cr} gets nearly 5 points gain over the TABLE 4

Experimental results on VQA-CP v2 test set and VQA v2 val set of state-of-the-art methods. **Best** and **second** performance are highlighted in each column. Methods with * use extra annotations (*e.g.*, human attention (HAT), explanations (VQA-X), or object label information). Methods with [†] use extra datasets at their pre-training stage. Methods with CGD are our reimplementation using released codes. Other results are reported in the original papers. The results are from 4 repeated experiments under different random seeds.

Mathad	Page		V	QA-CP tes	st			VQA	v2 val	
Method	Dase	All	Y/N	Num.	Others	CGD	All	Y/N	Num.	Others
GVQA [11]	-	31.30	57.99	13.68	22.14	-	48.24	72.03	31.17	34.65
UpDn [62]	-	$39.81{\scriptstyle \pm 0.05}$	$43.09{\scriptstyle\pm0.11}$	$12.01{\scriptstyle\pm0.09}$	$45.81{\scriptstyle \pm 0.01}$	$2.57{\scriptstyle\pm1.23}$	$63.57{\scriptstyle\pm0.22}$	$80.90{\scriptstyle \pm 0.10}$	$42.58{\scriptstyle \pm 0.27}$	$55.75{\scriptstyle\pm0.07}$
S-MRL [18]	-	38.46	42.85	12.81	43.20	-	63.10	-	-	-
HINT* [75]	UpDn	47.50	67.21	10.67	46.80	-	63.38	81.18	42.14	55.66
SCR* [76]	UpDn	49.45	72.36	10.93	48.02	-	62.2	78.8	41.6	54.4
AdvReg. [46]	UpDn	41.17	65.49	15.48	35.48	-	62.75	79.84	42.35	55.16
RUBi [18]	UpDn	$46.68{\scriptstyle \pm 0.89}$	$68.43{\scriptstyle \pm 3.86}$	$11.64{\scriptstyle \pm 0.26}$	$44.53{\scriptstyle\pm0.17}$	$7.88{\scriptstyle \pm 1.41}$	$58.74{\scriptstyle \pm 0.54}$	$67.17{\scriptstyle\pm 6.42}$	$39.85{\scriptstyle \pm 1.40}$	$54.12{\scriptstyle \pm 0.84}$
LM [17]	UpDn	$49.13{\scriptstyle \pm 1.07}$	$72.38{\scriptstyle\pm3.12}$	$14.49{\scriptstyle\pm0.71}$	$46.42{\scriptstyle \pm 0.04}$	$9.17{\scriptstyle \pm 2.10}$	$63.46{\scriptstyle \pm 0.17}$	$81.15{\scriptstyle \pm 0.04}$	$42.27{\scriptstyle\pm0.21}$	$55.61{\scriptstyle \pm 0.34}$
LMH [17]	UpDn	$53.30{\scriptstyle \pm 0.70}$	$73.47{\scriptstyle\pm0.38}$	$31.90{\scriptstyle\pm4.37}$	$47.92{\scriptstyle \pm 0.10}$	$10.54{\scriptstyle \pm 0.68}$	$58.06{\scriptstyle\pm2.18}$	$72.60{\scriptstyle \pm 6.73}$	$37.33{\scriptstyle \pm 1.75}$	$52.39{\scriptstyle\pm2.59}$
DLP [77]	UpDn	48.87	70.99	18.72	45.57	-	57.96	76.82	39.33	48.54
GVQE* [78]	UpDn	48.75	-	-	-	-	<u>64.04</u>	-	-	-
CSS* [22]	UpDn	$40.32{\scriptstyle \pm 0.59}$	$42.15{\scriptstyle\pm1.28}$	$12.44{\scriptstyle \pm 0.27}$	$46.76{\scriptstyle \pm 0.33}$	$8.58{\scriptstyle \pm 2.34}$	$62.34{\scriptstyle\pm1.85}$	$79.50{\scriptstyle\pm2.20}$	$42.11{\scriptstyle \pm 1.06}$	$55.50{\scriptstyle\pm1.32}$
CF-VQA(Sum) [79]	UpDn	53.69	<u>91.25</u>	12.80	45.23	-	63.65	<u>82.63</u>	44.01	54.38
RUBi [18]	S-MRL	47.11	68.65	20.28	43.18	-	61.16	-	-	-
GVQE* [78]	S-MRL	50.11	66.35	27.08	46.77	-	63.18	-	-	-
CF-VQA(Sum) [79]	S-MRL	54.95	90.56	21.88	45.36	-	60.76	81.11	43.48	49.58
MFE [51]	LMH	54.55	74.03	49.16	45.82	-	-	-	-	-
CSS* [22]	LMH	$58.27{\scriptstyle\pm0.05}$	$81.76{\scriptstyle \pm 1.34}$	$48.99{\scriptstyle\pm 5.85}$	$47.99{\scriptstyle \pm 0.12}$	$6.34{\scriptstyle \pm 1.75}$	$53.42{\scriptstyle \pm 0.26}$	$58.32{\scriptstyle\pm2.49}$	$37.99{\scriptstyle\pm1.15}$	$54.53{\scriptstyle\pm1.01}$
SAR [†] [80]	LMH	<u>62.51</u>	76.40	<u>59.40</u>	<u>56.09</u>	-	<u>65.79</u>	77.26	<u>52.71</u>	<u>60.52</u>
GGD_{gs}^{dq}	UpDn	$56.95{\scriptstyle \pm 0.34}$	$87.02{\scriptstyle\pm0.30}$	$25.97{\scriptstyle\pm1.35}$	$49.40{\scriptstyle\pm0.28}$	$\underline{15.24}_{\pm 0.93}$	$59.51{\scriptstyle\pm1.34}$	$74.77{\scriptstyle\pm3.22}$	$39.46{\scriptstyle\pm1.26}$	53.50±0.69
GGD^{dq}_{cr}	UpDn	$59.37{\scriptstyle\pm0.26}$	$88.23{\scriptstyle\pm0.29}$	$38.11{\scriptstyle \pm 1.05}$	$49.82{\scriptstyle \pm 0.40}$	$13.31{\scriptstyle \pm 1.69}$	$62.15{\scriptstyle \pm 0.93}$	$79.25{\scriptstyle\pm2.19}$	$42.43{\scriptstyle\pm0.21}$	$54.66{\scriptstyle \pm 0.32}$

BiDAF baseline that is already strong enough. This well demonstrates the power of the greedy learning in focusing on the hard/valuable samples from a biased dataset.

However, the limited improvement in the original setting is likely caused by Self-Ensemble, where the baseline model BiDAF itself can hardly capture useful biases from the dataset. If we can define a better biased model that can access the word-level similarity, we may achieve better performance without extra training data.

4.3 Visual Question Answering

4.3.1 Dataset

Data bias problems in multi-modal tasks are more challenging, where multiple data sources from different modalities should be jointly considered. In this section, we choose Visual Question Answering (VQA) as the representative multi-modal task for demonstration. Neural networks [82], [83], [84], [85], [86] that model the correlations between vision and language have shown remarkable results on large-scale benchmark datasets [9], [44], [87], [88], but most VQA methods tend to rely on existing idiosyncratic biases in the datasets [9], [10] and show poor generalization ability to out-of-domain data. In this section, we demonstrate the effectiveness of GGD on the challenging datasets VQA-CP v2 [11] and GQA-OOD [39]. VQA v2 [44] is a commonly used VQA dataset composed of real-world images from MSCOCO with the same train/validation/test splits. For each image, an average of three questions are generated, and 10 answers are collected for each image-question pair from human annotators. Following previous works, we take the answers that appeared more than 9 times in the training set as candidate answers, which produces 3129 answer candidates.

VQA-CP v2 [11] dataset is derived from the VQA 2.0 [44] but contains different answer distribution per question type between training and validation splits. Since it has different distribution on the train and test sets, the performance on this dataset better reflects models' generalization ability. VQA-CP v2 consists of 438,183 samples in the train set and 219,928 samples in the test set.

GQA-OOD [39] divides the test set of GQA [88] into majority (head) and minority (tail) groups based on the answer frequency within each 'local group', which is a unique combination of answer type (*e.g.*, colors) and the main concept (*e.g.*, 'bag', 'chair', *etc.*). The models are trained on the original GQA-balanced but tested on different finegrained local groups.

4.3.2 Experimental Setups

In the following experiments, we use UpDn [62] as our base model and the images are represented as object features

TABLE 5 Ablation study for different versions GGD on VQA-CP v2 test set and VQA v2 val set. ID indicates the overall Accuracy on VQA v2 val. **Best** results are highlighted in the columns. SE denotes the self-ensemble.

Method	All	Y/N	Others	Num.	CGD	ID
Baseline	39.89	43.01	45.80	11.88	3.91	63.79
SUM-DQ	35.46	42.66	38.01	12.38	3.10	56.85
LMH+RUBi	51.54	74.55	47.41	22.65	6.12	60.68
GGD_{gs}^d	48.27	70.75	47.53	13.42	14.31	62.79
GGD_{gs}^q	43.72	48.17	48.78	14.24	6.70	61.23
GGD_{gs}^{dq}	57.12	87.35	49.77	26.16	16.44	59.30
GGD^d_{cr}	50.93	78.50	47.30	12.92	10.28	62.17
GGD^q_{cr}	55.81	88.59	48.74	20.96	13.46	62.48
GGD_{cr}^{dq}	59.57	88.44	50.23	36.95	13.92	63.11
GGD^{se}_{gs}	44.53	50.98	48.90	18.24	6.08	59.30
GGD^{dse}_{gs}	56.33	86.43	49.32	24.37	14.47	61.03
GGD^{se}_{cr}	54.42	80.26	48.64	29.42	6.70	61.09
GGD_{cr}^{dse}	57.08	85.75	48.88	36.54	11.62	62.27

pre-extracted with Faster R-CNN [89]. The implementation details and the experiments on other base models are provided in the Appendix. All methods are measured with Accuracy and Correct Grounding Difference (CGD) proposed in [29]. CGD evaluates whether the visual information is well taken in answer decision.

For ablation studies, we present five different versions of GGD. \mathbf{GGD}^d only removes the distribution bias. \mathbf{GGD}^q only models the shortcut bias. \mathbf{GGD}^{dq} makes use of both the distribution bias and the shortcut bias. \mathbf{GGD}^{se} is the selfensemble version GGD, which takes the baseline model itself as the biased model. \mathbf{GGD}^{dse} removes the distribution bias before Self-Ensemble. The implementation details of above five ablations are provided in the Appendix.

4.3.3 Experimental Results

GGD can handle multiple biases. In the first group of ablation study, we compare with the other two ensemble strategies to verify the effectiveness of the greedy learning. **SUM-DQ** directly sums up the outputs of biased models and the base model. **LMH+RUBi** combines LMH [17] and RUBi [18]. It reduces distribution bias with LMH and shortcut bias with RUBi. The implementation details for these two ablations are provided in the Appendix.

As shown in TABLE 5, SUM-DQ performs even worse than vanilla UpDn. LMH+RUBi does not make use of both kinds of biases, whose Accuracy is just similar to that of LMH. On the other hand, both GGD_{gs} and GGD_{cr} surpass these two ablations by a large margin. This shows that the greedy strategy in GGD can really force the biased data to be learned with biased models in priority. As a result, the base model has to pay more attention to hard examples that are hard to solve based on the estimation of either distribution bias or shortcut bias. It needs to consider more visual information for the final decision.

In the second group of experiments, we directly compare GGD^d , GGD^q , and GGD^{dq} . As shown in TABLE 5, GGD_{gs}^{dq} surpasses single-bias versions GGD^d and GGD^q by ~10%. This well verifies that GGD can reduce multiple biases with



Fig. 5. Predicted distribution for three variants of GGD_{gs} .

the greedy learning procedure. The case analysis in Figure 5 provides a more qualitative evaluation. It shows that GGD^{*d*} uniforms predictions, which mainly improves Y/N as shown in TABLE 5. B_q works like "hard example mining" but will also introduce some noise (*e.g.*, "mirror" and "no" in this example) due to the unbalanced data distribution. GGD^{*d*q} can make use of both biases. Reducing B_d at first can further help the discovery of the hard examples with B_q and encourage the base model to capture essential visual information.

Implicitly De-bias with Self-Ensemble. In order to further discuss the generalizability of GGD, we also test a more flexible Self-Ensemble fashion GGD^{se} on VQA-CP. As shown in TABLE 5, GGD^{se} still surpasses UpDn without predefined biased features. Moreover, if we first remove distribution bias before Self-Ensemble, the performance of GGD^{dse} is comparable to existing state-of-the-art methods as well.

 GGD_{gs} vs. GGD_{cr} . As shown in TABLE. 5, GGD_{cr} largely alleviates the degradation on in-distribution test data VQA v2 val, which is even comparable to the original UpDn baseline. Moreover, It also gets better performance on VQA-CP under all GGD^d , GGD^q , GGD^{dq} , and GGD^{se} . The major improvement comes from the "Num." and "Other" question types which contain fewer samples in the training set. GGD_{gs} harms the performance on these question types because of the greedily discarding of samples that is easy to answer.

Comparison with State-of-the-art Methods. We compare our best performed model GGD^{dq}_{cr} with existing state-of-theart bias reduction techniques, including visual-groundingbased methods HINT [75], SCR [76], ensemble-based methods AdvReg. [46], RUBi [18], LM (LMH) [17], MFE [51], new-question-encoding-based methods GVQE [78], DLP [79], counterfactual-based methods CF-VQA [79], CSS [22], recent proposed regularization method MFE [51], and SAR [80] that models VQA as Visual Entailment with pre-trained model.

As shown in TABLE 4, GGD_{cr}^{dq} achieves state-of-the-art performance without extra bias annotation. It outperforms the baseline model UpDn by 20% higher in terms of Accuracy and 10% higher in terms of CGD, which verifies the effectiveness of GGD on both answer classification and visual-grounding ability.

For the comparison of question-type-wise results, incorporating GGD improves the performance for all the question types, especially the more challenging "other" question type [90]. CF-VQA [79] performs the best in Y/N, but worse than our methods in all the other question types and metrics.



What is the object in the water?

Fig. 6. Qualitative Evaluation for GGD_{gs}^{dq} . We provide a comparison between UpDn and GGD_{gs}^{dq} on the visualization of the most sensitive regions and confidence of the top-5 answers. Red answers denote the ground-truth.



Fig. 7. Comparison between GGD_{gs}^{dq} and GGD_{cr}^{dq} . The major improvements are reflected on counting problems and questions that rarely appear in the training data.



Fig. 8. Failure cases for GGD^{dq}_{cr}. From top to bottom, the failure cases are 1) counting problems; 2) Synonym answers; 3) incorrect visual evidences.

LMH [17], LMH-MFE [51], and LMH-CSS [22] work well on Num. questions. Comparing with LM and LMH, it is obvious that the performance gains in Num. are mainly due to the additional regularization for entropy. However, methods with entropy regularization drop nearly 10% on VQA v2. This indicates that these models may over-correct the bias. SAR [80] achieves the best performance on both VQA-CP v2 and VQA v2, using extra datasets for model pre-training. On the other hand, GGD_{cr} improves both Num. Accuracy and in-distribution performance only with the Curriculum Regularization and work well without any extra data sources. **Qualitative Evaluation**. Examples in Fig. 6 illustrate how GGD_{gs}^{dq} makes difference compared with the baseline UpDn [62]. The first example is about using visual information for inference. Despite offering the right answer "yes", the prediction from UpDn is not based on the right visual grounding result. In comparison, GGD_{gs} correctly grounds the giraffe that is eating leaves. The second example is a case of reducing language prior apart from Yes/No questions. UpDn answers "boat" just based on the language context "in the water", while GGD_{gs}^{dq} provides correct answers "tv" and "television" with more salient visual grounding. These examples qualitatively verify our improvement on both Accuracy and visual explanation for the predictions.

TABLE 6

Experimental results on GQA-OOD test-dev. g denotes the method addressing global group distribution bias, l is method addressing the local group distribution bias, and q is the method addressing the question shortcut bias. "Avg" is the mean accuracy of head and tail groups. Methods with only Avg are reported from [24].

Method	All Head		Tail	Avg
UpDn [62]	46.93±0.29	$49.41{\scriptstyle \pm 0.26}$	$42.73{\scriptstyle \pm 0.57}$	$46.07{\scriptstyle\pm 0.42}$
LfF [25]	47.06±0.23	$49.53{\scriptstyle \pm 0.63}$	$42.99{\scriptstyle \pm 0.47}$	$46.26{\scriptstyle \pm 0.09}$
SD [26]	47.59 ± 0.33	$50.05{\scriptstyle\pm0.32}$	$44.49{\scriptstyle\pm1.15}$	$47.27{\scriptstyle\pm0.42}$
RUBi ^q [18]	45.51±0.18	$47.87{\scriptstyle \pm 0.02}$	$41.67{\scriptstyle \pm 0.47}$	$44.71{\scriptstyle\pm 0.23}$
RUBi ^g [18]	$7.15{\scriptstyle \pm 0.53}$	$6.86{\scriptstyle \pm 0.91}$	6.60 ± 1.22	$6.573{\scriptstyle \pm 1.02}$
RUBi ^l [18]	20.48±3.84	$22.30{\scriptstyle \pm 4.04}$	17.51 ± 3.67	19.91±3.79
Up Wt ^g [5]	-	-	-	26.4
Up Wt ^{<i>l</i>} [5]	-	-	-	26.2
LNL ^g [20]	-	-	-	32.4
LNL ¹ [20]	-	-	-	10.7
GGD_{gs}^q	$47.41{\scriptstyle \pm 0.45}$	50.07±1.19	$43.09{\scriptstyle\pm1.06}$	$46.58{\scriptstyle\pm0.31}$
GGD_{gs}^g	48.96 ± 0.08	$52.07{\scriptstyle\pm0.12}$	$44.00{\scriptstyle\pm0.27}$	$48.03{\scriptstyle\pm0.13}$
GGD_{gs}^l	47.84 ± 0.49	$50.36{\scriptstyle \pm 0.60}$	$43.74{\scriptstyle \pm 0.33}$	$47.05{\scriptstyle\pm0.45}$
GGD_{gs}^{gq}	47.15±0.29	$49.62{\scriptstyle \pm 0.43}$	$43.27{\scriptstyle\pm0.89}$	$47.01{\scriptstyle\pm0.47}$
GGD^{lq}_{gs}	$48.25{\scriptstyle\pm0.53}$	50.74±0.99	$44.18{\scriptstyle \pm 0.24}$	$47.46{\scriptstyle \pm 0.39}$
GGD^q_{cr}	47.87±0.59	$50.19{\scriptstyle \pm 0.63}$	$44.25{\scriptstyle\pm0.31}$	$47.22{\scriptstyle \pm 0.47}$
GGD^g_{cr}	$48.09{\scriptstyle\pm0.35}$	$51.27{\scriptstyle\pm0.63}$	$43.25{\scriptstyle\pm0.35}$	$47.26{\scriptstyle \pm 0.44}$
GGD_{cr}^{l}	48.01 ± 0.60	$51.23{\scriptstyle \pm 0.87}$	$42.74{\scriptstyle \pm 0.17}$	$46.99{\scriptstyle \pm 0.51}$
GGD_{cr}^{gq}	$49.21{\scriptstyle \pm 0.08}$	$52.01{\scriptstyle \pm 0.30}$	$44.67{\scriptstyle\pm0.68}$	$48.34{\scriptstyle \pm 0.19}$
GGD_{cr}^{lq}	$47.03{\scriptstyle \pm 0.52}$	$49.37{\scriptstyle\pm0.71}$	$43.05{\scriptstyle\pm0.62}$	$46.21{\scriptstyle\pm0.42}$
GGD^{se}_{gs}	47.00 ± 0.10	$49.16{\scriptstyle \pm 1.44}$	42.61 ± 1.14	$4\overline{5.89}{\scriptstyle \pm 0.37}$
GGD^{se}_{cr}	$47.50{\scriptstyle \pm 0.35}$	51.10±0.39	$42.06{\scriptstyle \pm 0.75}$	$46.58{\scriptstyle \pm 0.41}$

Fig. 7 provides qualitative illustration for improvements achieved by GGD_{cr} . Compared with GGD_{gs} , GGD_{cr} majorly improves on the "Num." and "Other" question types. Questions about facial expression and counting do not frequently appear in the dataset. If these questions can be correctly predicted by fitting the biases, the base model in GGD_{gs} may not have enough data to learn a good representation.

Fig. 8 shows the examples of failure cases from GGD_{cr}^{dq} . The model appears to be weak on the complicated counting problem. Some failure cases are due to missing annotation in the dataset ("decoration" can also be regarded as the right answers). Although making wrong predictions, answers for failure cases in the last row are still consistent with visual explanations rather than language bias, which further indicates that GGD^{dq} truly makes use of the visual information.

4.3.4 Experimental Results on GQA-OOD

Biases. GQA-OOD is a more challenging dataset for Visual Question Answering. It has biases from multiple sources including imbalanced answer distribution, visual concept co-occurrences, question word correlations, and question type/answer distribution. Since the training set for GQA-OOD is the manually balanced GQA-balanced-train split, it is hard to specify the explicit biases to ensure that the models can generalize to even the rarest local groups. Following [24], apart from question shortcut bias similar to that in VQA-CP [11], we define two kinds of distribution bias according



Fig. 9. The loss ratio of the hard examples from baseline model and ${\rm GGD}_{cr}$ base model.

to global group labels (115 groups) and local group labels (133328 groups), and the corresponding models are denoted as GGD^{g} and GGD^{l} .

Comparison with State-of-the-art Methods. We compare GGD_{gs} with implicit de-bias methods LfF [25] and SD [26]; explicit de-bias methods RUBi [18], Up Wt [5], and LNL [20]. As shown in TABLE 6, all three previous explicit methods fail on both global and local group distribution bias, performing even worse than the original baseline UpDn. RUBi [18] with the question-only branch also degrades on GQA-OOD compared with the baseline. Implicit methods LfF [25] and SD [26] are more stable on handling complicated biases in GQA. SD [26] achieves the highest accuracy on the Tailed group. This indicates that both distribution bias and shortcut bias in GQA-OOD are not as obvious as those in VQA-CP, since the data from GQA is synthetic and the train split has been manually balanced.

On the other hand, GGD works well with hard-example mining. It surpasses the baseline under all bias settings and is comparable to existing implicit methods. If the biased models can not make a prediction with high confidence, the pseudo labels for the base model remain almost unchanged for most of the data. Although the biased models cannot well model the biases in the dataset, GGD will not harm the performance like previous explicit de-bias methods.

Ablation Study. According to the ablation studies, GGD^g addressing distribution bias on global groups works better compared with GGD^l on the local groups. Sequentially reducing the distribution bias and the shortcut bias will improve the Tail group Accuracy but slightly degrades the Head group accuracy. GGD_{cr}^{gq} achieves the highest overall accuracy because it alleviates the over-estimated bias in GGD_{gs}^{gg} . However, since the biased models do not capture biases with high confidence, GGD_{cr} does not show much difference compared with GGD_{gs} under most of the settings.

4.4 Discussion

Hard Example Mining Mechanism. To demonstrate the hard example mining mechanism of GGD, we provide analysis on the training process. We first evaluate whether the base model focuses on examples. We define the hard ratio R_h as

$$R_h = \frac{L_{hard}}{L_{all}},\tag{22}$$



Fig. 10. The original Labels and Pseudo Labels from GGD_{gs} . The left figure is the label changes in the Biased MNIST training set, while the right figure is that of "is this" question type in VQA-CP v2.

where L_{hard} is the loss of hard examples (*i.e.*, samples that do not choose corresponding background color), L_{all} is the loss for all samples. We calculate the accumulated loss for every 118 iterations of 4 repeated experiments with different random seeds. For more obvious comparison, experiments are done on GGD_{gs}^{1k} ($\lambda_t = 1$) on Biased MNIST with $\rho_{\text{train}} =$ 0.998. As shown in Fig. 9, although the proportion of hard examples is no more than 1%, R_h from the baseline is always over 0.6 and hardly decreases. The loss is trapped in local minimum due to the low *average* training error dominated by the majority groups of data. On the other hand, R_h from the base model of GGD is always lower than the baseline in any observed iteration and continuously decreases along with the training process. This reflects that GGD can well handle the hard examples compared with vanilla ResNet-18.

Bias Over-estimation. An interesting phenomenon is that GGD_{as} degrades on VQA v2 [44], but is relatively stable on the in-distribution test of Biased MNIST [12]. We find that this is due to the distribution bias on VQA v2. As shown in Fig. 10, taking "is this" question type as an example, GGD_{qs} will amplify the distribution bias and result in an "inverse biased distribution". On the other hand, the pseudo labels for Biased MNIST are still balanced because the texture bias is independent of the label distribution. This can also partially explain the improvement of GGD_{cr} on VQA, where the Curriculum Regularization also reduces such "bias over-estimation", apart from better low-level representation learning ability. In practice, one can get a more balanced classifier by selecting better λ_t according to the bias level of the dataset. It can also be a valuable research to adaptively estimate the bias level of a dataset in the future.

Limitations. Although the bias over-estimation problem in our previous GGE model has been alleviated with the Curriculum Regularization, there is still two major shortcomings of GGD. First, the hard-example mining mechanism in GGD is an instance-level sample re-weighting. If all samples are following a certain spurious correlations, GGD will fail to discover it as a spurious correlation (*e.g.*, $\rho_{\text{train}} = 1.0$ in Biased-MNIST). The gradients from the biased models will decline to 0. Even though the spurious feature is identified with the greedily learned biased models, the base model cannot learn a de-biased feature accordingly. Eq.12 indicates that we can also directly optimize $\log p(y|x^b)$ with the biased feature x^b , which can be obtained from the optimal biased model $h(x^b; \phi^*)$. We will investigate how to select network activations according to the bias models towards featurelevel ensemble in the future.

Second, if the biased model can well capture the biases in the dataset, GGD will largely improve both the indistribution and the out-of-distribution performance. However, when the biased model can not perfectly disentangle the spurious correlations, the improvement from GGD is limited (see experiments on Adversarial SQuAD and GQA-OOD). Although the Self-Ensemble fashion GGD^{se} can implicitly model the biases, it largely relies on the bias level of the dataset and the existence of hard examples in datasets. It can be a future work to design a more robust strategy that can capture spurious correlations needless of prior knowledge.

5 CONCLUSION

In this paper, we propose General Greedy De-bias Learning, a general de-bias framework with flexible regularization and wide applicability. Accompanied with Curriculum Regularization, the relaxed GGD_{cr} comes to a good trade-off between in-distribution and out-of-distribution performance. Experiments on image classification, Adversarial QA, and VQA demonstrate the effectiveness of GGD under both taskspecific biased models and self-ensemble fashion without prior knowledge on both ID and OOD scenarios.

In theory, the core of our method is the greedy strategy, which sequentially learns biased models in priority. One can also replace the regularization with better metrics that are able to measure the distance between the predictions and the labels. It may further improve the performance on specific tasks. In the future, we will try de-bias learning at the feature level and design a better strategy to capture spurious correlations needless of dataset-specific knowledge.

ACKNOWLEDGMENTS

This work was supported in part by the National Key R&D Program of China under Grant 2018AAA0102000, in part by National Natural Science Foundation of China: 62022083, 62236008 and 61931008, and in part by the Beijing Nova Program under Grant Z201100006820023.

REFERENCES

- Z. Qi, S. Wang, C. Su, L. Su, Q. Huang, and Q. Tian, "Selfregulated learning for egocentric video activity anticipation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [2] A. Torralba and A. Efros, "Unbiased look at dataset bias," in Proceedings of the 2011 IEEE Conference on Computer Vision and Pattern Recognition, 2011, pp. 1521–1528.
- [3] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," in *Proceedings of the International Conference on Learning Representations*, 2017.
- [4] R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann, "Shortcut learning in deep neural networks," *Nature Machine Intelligence*, vol. 2, no. 11, pp. 665–673, 2020.
- [5] S. Sagawa, A. Raghunathan, P. W. Koh, and P. Liang, "An investigation of why overparameterization exacerbates spurious correlations," in *International Conference on Machine Learning*. PMLR, 2020, pp. 8346–8356.

- [6] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, "Squad: 100,000+ questions for machine comprehension of text," in *Proceedings of the* 2016 Conference on Empirical Methods in Natural Language Processing, 2016, pp. 2383–2392.
- [7] R. Jia and P. Liang, "Adversarial examples for evaluating reading comprehension systems," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2017, pp. 2021– 2031.
- [8] S. Min, E. Wallace, S. Singh, M. Gardner, H. Hajishirzi, and L. Zettlemoyer, "Compositional questions do not necessitate multi-hop reasoning," in *Proceedings of the 57th Annual Meeting* of the Association for Computational Linguistics. Florence, Italy: Association for Computational Linguistics, Jul. 2019, pp. 4249–4257. [Online]. Available: https://www.aclweb.org/anthology/P19-1416
- [9] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh, "Making the v in vqa matter: Elevating the role of image understanding in visual question answering," in *Proceedings of the IEEE Conference* on Computer Vision and Pattern Recognition, 2017, pp. 6904–6913.
- [10] K. Kafle and C. Kanan, "An analysis of visual question answering algorithms," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1965–1973.
- [11] A. Agrawal, D. Batra, D. Parikh, and A. Kembhavi, "Don't just assume; look and answer: Overcoming priors for visual question answering," in *Proceedings of the IEEE Conference on Computer Vision* and Pattern Recognition, 2018, pp. 4971–4980.
- [12] H. Bahng, S. Chun, S. Yun, J. Choo, and S. J. Oh, "Learning debiased representations with biased representations," in *International Conference on Machine Learning*. PMLR, 2020, pp. 528–539.
- [13] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 3730–3738.
- [14] Y. He, Z. Shen, and P. Cui, "Towards non-iid image classification: A dataset and baselines," *Pattern Recognition*, vol. 110, p. 107383, 2021.
- [15] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma, "Learning imbalanced datasets with label-distribution-aware margin loss," in *Advances in Neural Information Processing Systems*, 2019.
- [16] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, "Class-balanced loss based on effective number of samples," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9268–9277.
- [17] C. Clark, M. Yatskar, and L. Zettlemoyer, "Don't take the easy way out: Ensemble based methods for avoiding known dataset biases," in Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, 2019, pp. 4060–4073.
- [18] R. Cadene, C. Dancette, M. Cord, D. Parikh et al., "Rubi: Reducing unimodal biases for visual question answering," in Advances in neural information processing systems, 2019, pp. 841–852.
- [19] S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang, "Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization," in *Proceedings of the International Conference on Learning Representations*, 2019.
- [20] B. Kim, H. Kim, K. Kim, S. Kim, and J. Kim, "Learning not to learn: Training deep neural networks with biased data," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9012–9020.
- [21] E. Tartaglione, C. A. Barbano, and M. Grangetto, "End: Entangling and disentangling deep representations for bias correction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 13508–13517.
- [22] L. Chen, X. Yan, J. Xiao, H. Zhang, S. Pu, and Y. Zhuang, "Counterfactual samples synthesizing for robust visual question answering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 800–10 809.
- [23] D. Teney, E. Abbasnedjad, and A. van den Hengel, "Learning what makes a difference from counterfactual examples and gradient supervision," in *Proceedings of the European conference on computer* vision, 2020, pp. 580–599.
- [24] R. Shrestha, K. Kafle, and C. Kanan, "An investigation of critical issues in bias mitigation techniques," arXiv preprint arXiv:2104.00170, 2021.
- [25] J. H. Nam, H. Cha, S. Ahn, J. Lee, and J. Shin, "Learning from failure: De-biasing classifier from biased classifier," in 34th Conference on Neural Information Processing Systems (NeurIPS) 2020, vol. 33, 2020, pp. 20673–20684.

- [26] M. Pezeshki, S.-O. Kaba, Y. Bengio, A. Courville, D. Precup, and G. Lajoie, "Gradient starvation: A learning proclivity in neural networks," arXiv preprint arXiv:2011.09468, 2020.
- [27] Z. Huang, H. Wang, E. P. Xing, and D. Huang, "Self-challenging improves cross-domain generalization," in *Proceedings of the European Conference on Computer Vision*, 2020.
- [28] E. Kim, J. Lee, and J. Choo, "Biaswap: Removing dataset bias with bias-tailored swapping augmentation," in *Proceedings of* the IEEE/CVF International Conference on Computer Vision (ICCV), October 2021, pp. 14992–15 001.
- [29] X. Han, S. Wang, C. Su, Q. Huang, and Q. Tian, "Greedy gradient ensemble for robust visual question answering," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1584–1593.
- [30] R. Xiong, Y. Chen, L. Pang, X. Cheng, Z.-M. Ma, and Y. Lan, "Uncertainty calibration for ensemble-based debiasing methods," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [31] Y. Niu and H. Zhang, "Introspective distillation for robust question answering," in Advances in Neural Information Processing Systems, vol. 34, 2021.
- [32] A. Kumar, T. Ma, P. Liang, and A. Raghunathan, "Calibrated ensembles can mitigate accuracy tradeoffs under distribution shift," in *The 38th Conference on Uncertainty in Artificial Intelligence*, 2022.
- [33] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proceedings of the 26th annual international conference on machine learning*, 2009, pp. 41–48.
- [34] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in 2009 IEEE conference on computer vision and pattern recognition. Ieee, 2009, pp. 248–255.
- [35] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [36] S. M. Xie, A. Kumar, R. Jones, F. Khani, T. Ma, and P. Liang, "Inn-out: Pre-training and self-training using auxiliary information for out-of-distribution robustness," in *International Conference on Learning Representations*, 2021.
- [37] A. Kumar, A. Raghunathan, R. M. Jones, T. Ma, and P. Liang, "Finetuning can distort pretrained features and underperform out-ofdistribution," in *International Conference on Learning Representations*, 2021.
- [38] R. Jia and P. Liang, "Adversarial examples for evaluating reading comprehension systems," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Copenhagen, Denmark: Association for Computational Linguistics, Sep. 2017, pp. 2021–2031. [Online]. Available: https://aclanthology.org/D17-1215
- [39] C. Kervadec, G. Antipov, M. Baccouche, and C. Wolf, "Roses are red, violets are blue... but should vqa expect them to?" in *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 2776–2785.
- [40] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, "Class-balanced loss based on effective number of samples," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 9268–9277.
- [41] S. Reddy, D. Chen, and C. D. Manning, "Coqa: A conversational question answering challenge," *Transactions of the Association for Computational Linguistics*, vol. 7, pp. 249–266, 2019.
- [42] E. Choi, H. He, M. Iyyer, M. Yatskar, W.-t. Yih, Y. Choi, P. Liang, and L. Zettlemoyer, "QuAC: Question answering in context," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Brussels, Belgium: Association for Computational Linguistics, Oct.-Nov. 2018, pp. 2174–2184. [Online]. Available: https://aclanthology.org/D18-1241
- [43] P. Zhang, Y. Goyal, D. Summers-Stay, D. Batra, and D. Parikh, "Yin and yang: Balancing and answering binary visual questions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5014–5022.
- [44] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick, and D. Parikh, "Vqa: Visual question answering," in *Proceed*ings of the IEEE international conference on computer vision, 2015, pp. 2425–2433.
- [45] Z. Shen, J. Liu, Y. He, X. Zhang, R. Xu, H. Yu, and P. Cui, "Towards out-of-distribution generalization: A survey," arXiv preprint arXiv:2108.13624, 2021.
- [46] S. Ramakrishnan, A. Agrawal, and S. Lee, "Overcoming language priors in visual question answering with adversarial regularization," in *Advances in Neural Information Processing Systems*, 2018, pp. 1541– 1551.

- [47] R. K. Mahabadi, Y. Belinkov, and J. Henderson, "End-to-end bias mitigation by modelling biases in corpora," in *Proceedings of the* 58th Annual Meeting of the Association for Computational Linguistics, 2020.
- [48] A. Gretton, O. Bousquet, A. Smola, and B. Schölkopf, "Measuring statistical dependence with hilbert-schmidt norms," in *International* conference on algorithmic learning theory. Springer, 2005, pp. 63–77.
- [49] P. P. Liang, C. Wu, L.-P. Morency, and R. Salakhutdinov, "Towards understanding and mitigating social biases in language models," in *International Conference on Machine Learning*. PMLR, 2021, pp. 6565–6576.
- [50] H. Wang, Z. He, Z. C. Lipton, and E. P. Xing, "Learning robust representations by projecting superficial statistics out," in *Proceedings* of the International Conference on Learning Representations, 2019.
- [51] I. Gat, I. Schwartz, A. Schwing, and T. Hazan, "Removing bias in multi-modal classifiers: Regularization by maximizing functional entropies," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [52] W. Zhu, H. Zheng, H. Liao, W. Li, and J. Luo, "Learning biasinvariant representation by cross-sample mutual information minimization," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 15002–15012.
- [53] Y. Hong and E. Yang, "Unbiased classification through biascontrastive and bias-balanced learning," Advances in Neural Information Processing Systems, vol. 34, 2021.
- [54] Z. Shen, P. Cui, T. Zhang, and K. Kunag, "Stable learning via sample reweighting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 5692–5699.
- [55] K. Kuang, R. Xiong, P. Cui, S. Athey, and B. Li, "Stable prediction with model misspecification and agnostic distribution shift," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 4485–4492.
- [56] V. Sanh, T. Wolf, Y. Belinkov, and A. M. Rush, "Learning from others' mistakes: Avoiding dataset biases without modeling them," in *Proceedings of the International Conference on Learning Representations*, 2021.
- [57] P. A. Utama, N. S. Moosavi, and I. Gurevych, "Towards debiasing nlu models from unknown biases," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing* (EMNLP), 2020.
- [58] C. Clark, M. Yatskar, and L. Zettlemoyer, "Learning to model and ignore dataset bias with mixed capacity ensembles," arXiv preprint arXiv:2011.03856, 2020.
- [59] X. Zhang, P. Cui, R. Xu, L. Zhou, Y. He, and Z. Shen, "Deep stable learning for out-of-distribution generalization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 5372–5382.
- [60] G. Blanchard, A. A. Deshmukh, U. Dogan, G. Lee, and C. Scott, "Domain generalization by marginal transfer learning," *Journal of Machine Learning Research*, vol. 22, pp. 1–55, 2021.
- [61] R. Christiansen, N. Pfister, M. E. Jakobsen, N. Gnecco, and J. Peters, "A causal framework for distribution generalization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [62] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, "Bottom-up and top-down attention for image captioning and visual question answering," in *Proceedings of the IEEE Conference* on Computer Vision and Pattern Recognition, 2018, pp. 6077–6086.
- [63] L. Mason, J. Baxter, P. L. Bartlett, and M. R. Frean, "Boosting algorithms as gradient descent," in Advances in neural information processing systems, 2000, pp. 512–518.
- [64] A. S. Rawat, A. K. Menon, W. Jitkrittum, S. Jayasumana, F. Yu, S. Reddi, and S. Kumar, "Disentangling sampling and labeling bias for learning in large-output spaces," in *Proceedings of the 38th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. Meila and T. Zhang, Eds., vol. 139. PMLR, 18–24 Jul 2021, pp. 8890–8901.
- [65] B. Kang, S. Xie, M. Rohrbach, Z. Yan, A. Gordo, J. Feng, and Y. Kalantidis, "Decoupling representation and classifier for longtailed recognition," *Proceedings of the International Conference on Learning Representations*, 2020.
- [66] S. Zhang, Z. Li, S. Yan, X. He, and J. Sun, "Distribution alignment: A unified framework for long-tail visual recognition," in *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 2361–2370.
- [67] M. Pagliardini, M. Jaggi, F. Fleuret, and S. P. Karimireddy, "Agree to disagree: Diversity through disagreement for better transferability," arXiv preprint arXiv:2202.04414, 2022.

- [68] Y. Freund, "Boosting a weak learning algorithm by majority," Information and computation, vol. 121, no. 2, pp. 256–285, 1995.
- [69] R. E. Schapire, "The strength of weak learnability," Machine learning, vol. 5, no. 2, pp. 197–227, 1990.
- [70] P. J. Bickel, Y. Ritov, and A. Zakai, "Some theory for generalized boosting algorithms," *Journal of Machine Learning Research*, vol. 7, no. May, pp. 705–732, 2006.
- [71] Y. Wald, A. Feder, D. Greenfeld, and U. Shalit, "On calibration and out-of-domain generalization," in *Advances in neural information* processing systems, vol. 34, 2021, pp. 2215–2227.
- [72] M. Yi, R. Wang, J. Sun, Z. Li, and Z.-M. Ma, "Improved ood generalization via conditional invariant regularizer," arXiv preprint arXiv:2207.06687, 2022.
- [73] D. Hendrycks, K. Lee, and M. Mazeika, "Using pre-training can improve model robustness and uncertainty," in *International Conference on Machine Learning*. PMLR, 2019, pp. 2712–2721.
- [74] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [75] R. R. Selvaraju, S. Lee, Y. Shen, H. Jin, S. Ghosh, L. Heck, D. Batra, and D. Parikh, "Taking a hint: Leveraging explanations to make vision and language models more grounded," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 2591– 2600.
- [76] J. Wu and R. Mooney, "Self-critical reasoning for robust visual question answering," in Advances in Neural Information Processing Systems, 2019, pp. 8604–8614.
- [77] C. Jing, Y. Wu, X. Zhang, Y. Jia, and Q. Wu, "Overcoming language priors in vqa via decomposed linguistic representations." in AAAI, 2020, pp. 11 181–11 188.
- [78] G. KV and A. Mittal, "Reducing language biases in visual question answering with visually-grounded question encoder," arXiv preprint arXiv:2007.06198, 2020.
- [79] Y. Niu, K. Tang, H. Zhang, Z. Lu, X.-S. Hua, and J.-R. Wen, "Counterfactual vqa: A cause-effect look at language bias," arXiv preprint arXiv:2006.04315, 2020.
- [80] Q. Si, Z. Lin, M. Zheng, P. Fu, and W. Wang, "Check it again: Progressive visual question answering via visual entailment," arXiv preprint arXiv:2106.04605, 2021.
- [81] M. Seo, A. Kembhavi, A. Farhadi, and H. Hajishirzi, "Bidirectional attention flow for machine comprehension," in *Proceedings of the International Conference on Learning Representations*, 2017.
- [82] P. Gao, H. You, Z. Zhang, X. Wang, and H. Li, "Multi-modality latent interaction network for visual question answering," arXiv preprint arXiv:1908.04289, 2019.
- [83] R. Hu, J. Andreas, M. Rohrbach, T. Darrell, and K. Saenko, "Learning to reason: End-to-end module networks for visual question answering," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017.
- [84] R. Hu, A. Rohrbach, T. Darrell, and K. Saenko, "Languageconditioned graph networks for relational reasoning," in *Proceedings* of the IEEE International Conference on Computer Vision, 2019, pp. 10294–10303.
- [85] J. Mao, C. Gan, P. Kohli, J. B. Tenenbaum, and J. Wu, "The neurosymbolic concept learner: Interpreting scenes, words, and sentences from natural supervision," arXiv preprint arXiv:1904.12584, 2019.
- [86] X. Han, S. Wang, C. Su, W. Zhang, Q. Huang, and Q. Tian, "Interpretable visual reasoning via probabilistic formulation under natural supervision," in *European Conference on Computer Vision*. Springer, 2020, pp. 553–570.
- [87] J. Johnson, B. Hariharan, L. van der Maaten, L. Fei-Fei, C. Lawrence Zitnick, and R. Girshick, "Clevr: A diagnostic dataset for compositional language and elementary visual reasoning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2901–2910.
- [88] D. A. Hudson and C. D. Manning, "Gqa: A new dataset for realworld visual reasoning and compositional question answering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6700–6709.
- [89] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [90] D. Teney, E. Abbasnejad, K. Kafle, R. Shrestha, C. Kanan, and A. Van Den Hengel, "On the value of out-of-distribution testing: An example of goodhart's law," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 407–417.

- [91] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in Interna*tional conference on machine learning.* PMLR, 2015, pp. 448–456. [92] J.-H. Kim, J. Jun, and B.-T. Zhang, "Bilinear attention networks," in
- Advances in Neural Information Processing Systems, 2018, pp. 1564-1574.
- [93] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in Proceedings of the 2014 conference on empirical methods in natural language processing, 2014, pp. 1532-1543.
- [94] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga et al., "Pytorch: An imperative style, high-performance deep learning library," in Advances in neural information processing systems, 2019, pp. 8026-8037
- [95] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in Proceedings of the IEEE international conference on computer vision, 2017, pp. 618-626.



Qingming Huang received the B.S. degree in computer science and Ph.D. degree in computer engineering from the Harbin Institute of Technology, Harbin, China, in 1988 and 1994, respectively. He is currently a Chair Professor with the School of Computer Science and Technology, University of Chinese Academy of Sciences. He has published over 500 academic papers in international journals, such as IEEE Transactions on Pattern Analysis and Machine Intelligence, IEEE Transactions on Image Processing, IEEE

Transactions on Multimedia, IEEE Transactions on Circuits and Systems for Video Technology, and top level international conferences, including the ACM Multimedia, ICCV, CVPR, ECCV, VLDB, and IJCAI. He was the Associate Editor of IEEE Transactions on Circuits and Systems for Video Technology and the Associate Editor of Acta Automatica Sinica. His research interests include multimedia computing, image/video processing, pattern recognition, and computer vision.



Xinzhe Han received the B.S. degree from Xidian University in 2017. He is currently pursuing the Ph.D. degree in the School of Computer Science and Technology, University of Chinese Academy of Sciences. His current research interests include visual question answering and trustable machine learning.



tronics engineering from Tsinghua University, Beijing, China, in 2006, and the Ph.D. degree from the Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China, in 2012. He is currently a Full Professor with the Key Laboratory of Intelligent Information Processing (CAS), Institute of Computing Technology, Chinese Academy of Sciences. He is also with Pengcheng Laboratory, Shenzhen. His research interests include image/video understand-

ing/retrieval, cross-media analysis and visual-textual knowledge extraction.



Qi Tian is currently a Chief Scientist in Artificial Intelligence at Cloud BU, Huawei. From 2018-2020, he was the Chief Scientist in Computer Vision at Huawei Noah's Ark Lab. He was also a Full Professor in the Department of Computer Science, the University of Texas at San Antonio (UTSA) from 2002 to 2019. During 2008-2009, he took one-year Faculty Leave at Microsoft Research Asia (MSRA). Dr. Tian received his Ph.D. in ECE from University of Illinois at Urbana-Champaign (UIUC) and received his B.E. in

Electronic Engineering from Tsinghua University and M.S. in ECE from Drexel University, respectively. Dr. Tian's research interests include computer vision, multimedia information retrieval and machine learning and published 600+ refereed journal and conference papers. His Google citation is over 44700+ with H-index 97. He was the co-author of best papers including IEEE ICME 2019, ACM CIKM 2018, ACM ICMR 2015, PCM 2013, MMM 2013, ACM ICIMCS 2012, a Top 10% Paper Award in MMSP 2011, a Student Contest Paper in ICASSP 2006, and coauthor of a Best Paper/Student Paper Candidate in ACM Multimedia 2019, ICME 2015 and PCM 2007. Dr. Tian research projects are funded by ARO, NSF, DHS, Google, FXPAL, NEC, SALSI, CIAS, Akiira Media Systems, HP, Blippar and UTSA. He received 2017 UTSA President's Distinguished Award for Research Achievement, 2016 UTSA Innovation Award, 2014 Research Achievement Awards from College of Science, UTSA, 2010 Google Faculty Award, and 2010 ACM Service Award. He is the associate editor of IEEE TMM, IEEE TCSVT, ACM TOMM, MMSJ, and in the Editorial Board of Journal of Multimedia (JMM) and Journal of MVA. Dr. Tian is the Guest Editor of IEEE TMM, Journal of CVIU, etc. Dr. Tian is a Fellow of IEEE.



Chi Su is currently a General Manager of Smart-More, Beijing. He received the PhD degree in the Institute of Digital Media, EECS, Peking University. His research include computer vision and machine learning, with focus on object detection, object tracking, and human identification and recognition.

APPENDIX A PROBABILISTIC JUSTIFICATION

We consider the distribution p(c|x):

$$p(y|x) = p(y|x^{b}, x^{-b})$$

$$\propto p(x^{-b}|y, x^{b})p(y|x^{b}) \triangleright \text{Bayes Rule}$$

$$= p(x^{-b}|y)p(y|x^{b}) \triangleright \text{Conditionally Independent}$$

$$= p(y|x^{b})\frac{p(y|x^{-b})p(x^{-b})}{p(y)} \triangleright \text{Bayes Rule}$$

$$\propto \frac{p(y|x^{b})}{p(y)}p(y|x^{-b})$$
(23)

Rearranging the log-likelihood of Eq.23 will lead to the probabilistic justification in Section 3.4.2. p(y|x) is more likely to align with $p(y|x^{-b})$ in bias-conflicting samples. Moreover, it shows that the distribution p(y) also has effects on the optimization of $p(y|x^{-b})$. This can partially explain the influence from the distribution bias.

APPENDIX B IMPLEMENTATION DETAILS

B.1 SimpleNet-1k

SimpleNet-1k is a fully convolutional CNN proposed in [12]. It contains four convolutional layers with 1×1 kernels and output channel {16, 32, 64, 128}. Each convolutional layer is followed by batch normalisation [91] and ReLU. The classification layer is consist of a Global Average Pooling (GAP) layer following a (128×10) linear projection.

All models for Biased MNIST is trained with batch size of 256 and Adam optimizer. The initial learning rate is set to be 1e-3.

B.2 Different versions of GGD for VQA

The optimization paradigm for GGD^d, GGD^q, GGD^q, GGD^{dq}, GGD^{se} and GGD^{dse} are shown in Fig. 11. V, Q and \tilde{A} denote images, questions, and answer predictions respectively. A is the human-annotated labels. $B_d : \{\hat{y}_i^d\}_{i=1}^N, B_q : \{\hat{y}_i^q\}_{i=1}^N$ and $B_{se} : \{\hat{y}_i^{se}\}_{i=1}^N$ indicate the prediction from distribution bias, question shortcut bias and self-ensemble bias respectively.

GGD^{*d*} only models distribution bias for ensemble. We define the distribution bias as answer distribution in the train set conditioned on question types:

$$\hat{y}_i^d = p(a_i|t_i),\tag{24}$$

where t_i denotes the type of question q_i . The reason for counting samples conditioned on question types is to maintain type information when reducing distribution bias. Question type information can only be obtained from the questions rather than the images, which does not belong to the language bias to be reduced.

The regularization for the base model is

$$L = \mathcal{L}\left(\tilde{A}, A\right) - \lambda_t \mathcal{L}_{CE}\left(\tilde{A}, B_d\right), \qquad (25)$$

where *A* is the predictions, and *A* is the labelled answers.

GGD^{*q*} only uses a question-only branch for shortcut bias. The shortcut bias is the semantic correlation between specific QA pairs. Similar to [18], we compose the question shortcut bias as a question-only branch

$$\hat{y}_i^q = c_q \left(e_q(q_i) \right), \tag{26}$$

where $c_q: Q \to \mathbb{R}^C$.

We first optimize the question-only branch with labelled answers

$$L_1 = \mathcal{L}(B_q, A). \tag{27}$$

The loss for base model is

$$L_2 = \mathcal{L}\left(\tilde{A}, A\right) - \lambda_t \mathcal{L}_{CE}\left(\tilde{A}, B_q\right).$$
(28)

 \mathbf{GGD}^{dq} uses both distribution bias and question shortcut bias. The loss for B_q is

$$L_1 = \mathcal{L}(B_q, A) - \lambda_t \mathcal{L}_{CE}(B_q, B_d).$$
(29)

The loss for base model is

$$L_2 = \mathcal{L}\left(\tilde{A}, A\right) - \lambda_t \mathcal{L}_{CE}\left(\tilde{A}, B_q + B_d\right).$$
(30)

 L_1 and L_2 are optimized iteratively.

GGD^{*se*} takes the joint representation $r_i = m(e_v(v_i), e_q(q_i))$ itself as the biased feature instead of predefined question-only branch in GGD^{*q*}, the biased prediction is

$$\hat{y}_i^{se} = c_s\left(r_i\right),\tag{31}$$

where $c_s : r \to \mathbb{R}^C$ is the classifier of the biased model.

We first optimize a baseline model with labelled answers

$$L_1 = \mathcal{L}(B_{se}, A), \tag{32}$$

The loss for base model is

$$L_2 = \mathcal{L}\left(\tilde{A}, A\right) - \lambda_t \mathcal{L}_{CE}\left(\tilde{A}, B_{se}\right).$$
(33)

 GGD^{dse} removes the distribution bias before Self-Ensemble, which is similar to GGD^{dq}

$$L_1 = \mathcal{L}\left(\sigma(B_{se}), A\right) - \lambda_t \mathcal{L}_{CE}\left(\sigma(B_{se}), B_d\right).$$
(34)

The loss for base model is

$$L_2 = \mathcal{L}\left(\sigma(\tilde{A}), A\right) - \lambda_t \mathcal{L}_{CE}\left(\sigma(\tilde{A}), \sigma(B_{se}) + B_d\right).$$
(35)

 L_1 and L_2 are optimized iteratively.

B.3 SUMB-DQ and LMH+RUBi

SUM-DQ directly sums up the outputs of biased models and the base model without greedy learning. The loss for the whole model is

$$L = \mathcal{L}(B_d + B_q + A, A), \tag{36}$$

where B_d is the predicted distribution bias, B_q is the predicted shortcut bias, \tilde{A} is the predictions and A is the labelled answers.

LMH+RUBi combines LMH [17] and RUBi [18]. It reduces distribution bias with LMH and shortcut bias with RUBi. The loss for RUBi is written as

$$L_{rubi}(\tilde{A}, A) = \mathcal{L}(\tilde{A} \odot \sigma(G_q), A) + \mathcal{L}(c_q(G_q), A), \quad (37)$$

where $G_q = g(e_q(q_i)), g(.) : Q \to \mathbb{R}^C$. Combining with LMH, the prediction is composed as

$$F(\tilde{A}, B, M) = \log \tilde{A} + h(M) \log B,$$
(38)



Fig. 11. Different versions of GGD for VQA. V, Q and \tilde{A} denote image, question, and answer prediction respectively. A is the human-annotated labels. B_d and B_q indicate the prediction from distribution bias and question shortcut bias respectively.

TABLE 7 Ablations for base model SimpleNet-7k on Biased MNIST.

$ ho_{ m train}/ ho_{ m test}$	0.990				0.995			0.999		
	0	0.1	0.990	0	0.1	0.995	0	0.1	0.999	
SimpleNet-7k	83.52	85.33	99.71	61.03	54.31	99.51	1.01	9.960	99.86	
ReBias [12]	86.39	88.15	99.81	78.17	81.32	99.86	25.17	33.58	99.88	
RUBi [18]	88.91	90.13	99.79	74.67	76.52	97.72	19.78	26.53	93.05	
GGD_{gs}	94.24	94.88	98.84	87.08	88.32	97.48	57.45	60.79	93.66	
GGD_{cr}	93.28	94.26	99.79	79.95	81.09	99.03	42.92	48.73	99.81	

where *M* and *B* are the fused feature and the bias in LMH, $h(.): M \to \mathbb{R}^C$. The combined loss function is

$$L = L_{rubi}(F(A, B, M), A) + wH(h(M)\log B),$$
 (39)

where H(.) is the entropy and w is a hyper-parameter.

B.4 UpDn

We use the publicly available reimplementation of $UpDn^1$ [62] for our baseline architecture, data preprocess and optimization in the VQA task.

Image Encoder. Following the popular bottom-up attention mechanism [62], we use a Faster R-CNN [89] based framework to extract visual features. We select the top-36 region proposals for each image $\mathbf{v} \in \mathbb{R}^{36 \times 2048}$.

Question Encoder. Each word is first initialized by 300dim GloVe word embeddings [93], then fed into a GRU with 1024-d hidden vector. The question representation is the last state of GRU $h_T \in \mathbb{R}^{1024}$.

1. https://github.com/hengyuan-hu/bottom-up-attention-vqa

Multi-modal Fusion. We use traditional linear attention between h_T and **v** for visual representation. and the final representation for classification is the Hadamard product of vision and question representation.

Question-only Classifier. The question-only classifier is implemented as two fully-connected layers with ReLU activations. The input question representation is shared with that in VQA base model.

Question types. We use 65 question types annotated in VQA v2 and VQA-CP, according to the first few words of the question (e.g., "What color is"). To save the training time, we simply use statistic answer distribution conditioned by question type in the train set as the prediction of distribution bias.

Optimization. Following UpDn [62], all the experiments are conducted with the Adamax optimizer for 20 epochs with learning rate initialized as 0.001. We train all models on a single RTX 3090 GUP with PyTorch 1.7 [94] and batch size 512.

TABLE 8 Ablations of base model BAN and S-MRL for VQA-CP v2.

Method	VQA-CP test									
Method	All	Y/N	Num.	Others	↑CGR	↓CGW	↑CGD			
S-MRL [18]	37.90	43.68	12.04	41.97	41.94	27.32	14.62			
$+ \text{GGD}_{gs}^{dq}$	54.03	79.66	20.77	46.72	38.10	22.42	15.68			
$+ GGD_{cr}^{dq}$	54.46	86.43	14.16	47.16	39.24	25.69	14.55			
BAN [92]	35.94	40.39	12.24	40.51	5.33	5.19	0.14			
$+ \text{GGD}_{gs}^{dq}$	50.75	74.56	20.59	46.54	20.87	16.85	4.98			
$+ \text{GGD}_{cr}^{dq}$	51.72	77.58	23.70	46.11	33.93	22.92	11.01			



Fig. 12. Saliency maps of the ResNet-18, bias model SimpleNet-1k, the base model from GGD_{gs} , and the base model from GGD_{cr} . The redder the pixel is, the more contributions it makes to prediction.

Data Preprocessing. Following previous works, we filter the answers that appear less than 9 times in the train set. For each instance with 10 annotated answers, we set the scores for labels that appear 1/2/3 times as 0.3/0.6/0.9, more than 3 times as 1.0.

APPENDIX C ABLATIONS OF BASE MODEL

GGD is agnostic for choice of the base model. In this section we provide extra experiments on Biased MNIST [12] and VQA-CP v2 [11].

For Biased MNIST, we do experiments on SimpleNet-7k following [12]. SimpleNet-7k has the same architecture with SimpleNet-1k introduced in Section B.1 but with convolution kernel size 7×7 . SimpleNet-1k is chosen for the biased model for all experiments in TABLE 7.

For VQA-CP, We do experiments on other base models BAN [92] and S-MRL [18]. The models are re-implemented based on officially released codes. For BAN, we set the number of Bilinear Attention blocks as 3. We choose the last bi-linear attention map of BAN and sum up along the question axis, which is referred to as the object attention for CGR and CGW. Although Accuracy of our reproduced S-MRL is a litter lower than that in [18], GGD^{dq} can improve the Accuracy over 10% and surpass most of the existing methods. As shown in the TABLE 8, GGD is a model-agnostic de-bias method, which can improve all three base models UpDn [62], S-MRL [18] and BAN [92] by a large margin.

APPENDIX D VISUALIZATION ON BIASED-MNIST

In this section we provide the saliency map visualizations of vanilla ResNet-18, SimpleNet-1k, and the GGD base model with grad-CAM [95]. All experiments are demonstrate on Biased-MNIST with $\rho_{\text{train}} = 0.999$.

As shown in Fig. 12, the vanilla ResNet-18 mainly focus on the background parts, since the background colour is highly correlated with the labels in the train stage. SimpleNet-1k provides an uniform saliency map due to the small perceptive field. In contrast, the visualization of GGD focus on the middle of the image, which means it capture more information about the digit number.

TABLE 9 Top-1 Accuracy of ResNet-32 on CIFAR-10-LT under different imbalance settings.

Imbalance Factor	1	0.2	0.1	0.02	0.01	Mean	$\uparrow \overline{\Delta}$	$\overline{\Delta}\%$
ResNet-32	92.03	86.28	83.30	79.83	66.61	81.61	-	-
	92.48 92.66	88.78 88.77	87.13 87.39	83.26 81.98	71.36 70.77	84.60 84.31	2.99 2.70	3.67% 3.31%
$\begin{array}{c} \text{GGD}_{gs}^{se} \\ \text{GGD}_{cr}^{se} \end{array}$	76.61 92.00	59.85 87.19	60.26 84.81	54.53 79.51	44.68 68.37	59.19 82.38	-22.42 0.77	-27.48% 0.94%

TABLE 10 Top-1 Accuracy of ResNet-32 on CIFAR-100-LT under different imbalance settings.

Imbalance Factor	1	0.2	0.1	0.02	0.01	Mean	$\uparrow \overline{\Delta}$	$\uparrow \overline{\Delta}\%$
ResNet-32	66.91	51.71	42.88	31.74	27.85	44.22	-	-
	67.55 67.13	59.44 54.77	49.03 45.52	35.88 34.85	31.09 31.03	48.60 46.66	4.38 2.44	9.91% 5.52%
$\begin{array}{c} & \\ & \text{GGD}_{gs}^{se} \\ & \text{GGD}_{cr}^{se} \end{array}$	48.62 66.04	34.50 52.16	32.70 44.48	27.69 32.35	21.47 31.81	33.00 45.37	-11.22 1.15	-25.38% 2.60%



Fig. 13. The gradient similarity between baseline model and ${\rm GGD}_{\it cr}$ base model.

APPENDIX E GRADIENT COMPARISON

To further compare the training process between GGD and the baseline, we provide the cosine similarity of the direction of the gradient for the feature before the classifier:

$$s = \frac{\mathbf{g}_{base} \mathbf{g}_{ggd}}{||\mathbf{g}_{base}||||\mathbf{g}_{ggd}||},\tag{40}$$

where ||.|| indicates L_2 norm of a vector. The gradients are also accumulated every 118 iterations and averaged with four repeated experiments with different random seeds. To guarantee that the baseline and the GGD base model are given the same initialization and trained with the same mini-batch of data at each iteration, this experiment is conducted on the Self-Ensemble GGD_{gs}^{se} . We plot the cosine similarity versus training iteration in Fig.13. It shows that the optimization direction with the same training data is extremely different between GGD and baseline, *s* is no more than 0.15 during training. This indicates that GGD learns a different feature compared with the baseline, which can better classify the hard examples in spite of the spurious correlations.

APPENDIX F

EXPERIMENTS ON LONG-TAILED IMAGE CLASSIFI-CATION

Dataset. The original CIFAR-10 (CIFAR-100) dataset contains 50,000 training images and 10,000 test images of size 32x32 uniformly falling into 10 (100) classes [35]. Cui et al. [40] created long-tailed versions by randomly removing training examples. In particular, $\mu = n_t/n_h$ controls the imbalance factor of the dataset, where n_h is the number of examples in the head class and n_t is the number of examples from the tailed class. By varying μ , we arrive at 5 training sets, respectively, with the imbalance factors of 1, 0.2, 0.1, 0.02, and 0.01, where $\mu = 1$ corresponds to the original datasets.

Biased Models. We test two settings of biased models in the experiments. The first biased model directly uses the statistical distribution of the training set as biased predictions, noted as GGD^d . The second is the self-ensemble GGD^{se}

Experimental Setups. We choose ResNet-32 as our baseline model. All experiments are conducted with the Adam optimizer for 250 epochs and the learning rate is initialized as 0.01.

Experimental Results. With the prior knowledge of the unbalanced distribution, both GGD_{gs} and GGD_{cr} can promisingly improve the performance of the long-tailed training data. However, in self-ensemble version without prior knowledge, GGD_{gs} will fail to estimate the biases of data. The baseline model itself cannot correctly reflect the data distribution. On the other hand, GGD_{cr} is much more robust. It at least keep the performance of the base model at all imbalanced levels without prior knowledge, even when the imbalanced factor $\mu = 1$.