051

052

053

054

Do Students Debias Like Teachers? On the Distillability of Bias Mitigation Methods

Anonymous Authors¹

Abstract

Knowledge distillation (KD) is an effective method for model compression and transferring knowledge between models. However, its effect on model's robustness against spurious correlations, shortcuts and task-irrelevant features that degrade performance on out-of-distribution data remains underexplored. This study investigates the effect of knowledge distillation on natural language inference (NLI) and image classification tasks, with a focus on the transferability of "debiasing" capabilities from teacher models to student models. Through extensive experiments, we illustrate several key findings: (i) the effect of KD on debiasing performance depends on the underlying debiasing method, the relative scale of the models involved, and the size of the training set; (ii) KD effectively transfers debiasing capabilities when teacher and student are similar in scale (number of parameters); (iii) KD may amplify the student model's reliance on spurious features, and this effect does not diminish as the teacher model scales up; and (iv) although the overall robustness of a model may remain stable post-distillation, significant variations can occur across different types of biases; and Given the above findings, we propose three effective solutions to improve the distillability of debiasing methods: developing high quality data for augmentation, implementing iterative knowledge distillation, and initializing student models with weights obtained from teacher models.

1. Introduction

Machine learning models are susceptible to biases or spurious correlations in datasets, commonly known to as "shortcuts" or "dataset biases". Models that rely on shortcuts can achieve high performance on in-domain test sets or overrepresented groups by exploiting superficial correlations between features and labels. However, these models suffer significant performance degradation on out of distribution or challenging test data, such as swapped subject and object (McCoy et al., 2019) in natural language understanding, or under-represented groups, such as "Male" subjects with "Blond Hair" (Liu et al., 2015) in image classification.

Despite recent advancements in bias mitigation (Guo et al., 2023; Chew et al., 2024; Li et al., 2023; Cheng & Amiri, 2024) and knowledge distillation (Stanton et al., 2021; Sultan, 2023), their integration is largely unexplored. This work studies the following research questions (RQs):

- **RQ1**: To what extent can knowledge distillation transfer debiasing capabilities between models?
- **RQ2**: Can knowledge distillation train less biased models compared to standard training?
- **RQ3**: Do different debiasing methods show consistent patterns in task performance and debiasing effective-ness before and after knowledge distillation?

Answering these questions will help us understand the efficacy of knowledge distillation in handling dataset biases, its underlying mechanisms, and its role in developing new training methods for bias mitigation.

We answer these questions by designing and conducting an empirical analysis on natural language understanding and image classification tasks. Our analyses show that: (i) the effect of knowledge distillation on debiasing performance depends on the underlying debiasing method, the relative scale of the models involved, and the size of the training set; (ii) knowledge distillation effectively transfers debiasing capabilities when teacher and student are similar in scale (number of parameters); (iii) knowledge distillation may amplify the student model's reliance on spurious features, and this effect does not diminish as the teacher model scales up; and (iv) although the overall robustness of a model may remain stable post-distillation, significant variations can occur across different types of biases; and (v) consistent

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

transfer patterns sometimes emerge, such as performance gap between teacher and student on out-of-distribution 057 (OOD) data, suggesting the possibility of predictable 058 changes in robustness after distillation. Given the above 059 findings, we propose three effective solutions to improve the 060 distillability of debiasing methods: developing high quality 061 data for augmentation, implementing iterative knowledge 062 distillation, and initializing student models with weights 063 obtained from teacher models.

The contributions of this paper are:

064

065

066

067

068

069

070

074

075

076

077

078 079 080

081

082

083

085

087

088

089

090

091

092

093

094

095

096

097

098

099

100

104

105

106

109

- to the best of our knowledge, we present the first study to investigate the effect of knowledge distillation on model robustness against dataset biases, and analyze the distillability of bias mitigation methods across both language and vision tasks;
 - our analysis reveals unique patterns in how knowledge distillation affects robustness to spurious correlations across different backbones and debiasing methods; and
 - we propose three strategies to improve the distillability of debiasing methods and provide insights for future development of bias mitigation techniques.

2. Knowledge Distillation and Debiasing

2.1. Problem Formulation

We investigate the effect of knowledge distillation (KD) on debiasing methods. We define *distillability of debiasing methods* as the amount of performance maintained before and after distilling a debiased model. We define *contribution of KD* as the performance improvement gained by training a debiasing method with KD over training without KD.

2.2. Notation and Training Setup

Let f and g denote models trained without knowledge distillation and with knowledge distillation respectively. In this paper, we use subscript \mathcal{T} and \mathcal{S} to denote teacher and student scales respectively. As illustrated in Figure 1, we train the following models for each debiasing method M_i : (i) we train M_i from scratch for both teacher and student scales to obtain $f_{\mathcal{T}}$ and $f_{\mathcal{S}}$, see Figure 1(a). (ii) Then for every scale $\mathcal{T} > \mathcal{S}$, we distill the knowledge from $f_{\mathcal{T}}$ to $g_{\mathcal{T}->\mathcal{S}}$, see Figure 1(b). Given a debiasing method M and the three models obtained above $(f_{\mathcal{T}}, f_{\mathcal{S}}, \text{ and } g_{\mathcal{T}->\mathcal{S}})$, we conduct the following comparisons:

• C1: Teacher (f_T) vs. Student $(g_{T->S})$. This comparison reveals if knowledge distillation can distill debiasing capability between models and if it affects model's robustness to spurious correlations, which answers RQ1 (§4).

- C2: Non-KD vs. KD, realized by comparing f_S vs. $g_{T->S}$. This comparison demonstrates if training bias mitigation networks can benefit from external knowledge from teacher models, which answers RQ2 (§5).
- C3: Comparison between debiasing methods (M_i vs. M_j). This comparison provides insights into differences between debiasing methods, which answers RQ3 (§6).

We note that when $\mathcal{T} = S$, C1 and C2 are essentially the same comparison. To avoid duplicate discussion, we will present results when $\mathcal{T} = S$ in C2.

3. Experimental Setup

For consistency and fair comparison with previous debiasing works in NLU (Jeon et al., 2023; Reif & Schwartz, 2023) and image classification (Kirichenko et al., 2023; LaBonte et al., 2023; Li et al., 2023), we adopt commonly used experimental setups, including choice of backbone models, datasets, evaluation protocols, and debiasing methods. In addition, all experiments are repeated three times with different random seeds to account for any stochastic effect.

Backbones We conduct experiments on a series of BERT (Devlin et al., 2019; Turc et al., 2019), T5 (Tay et al., 2022), ResNet (He et al., 2016), and ViT (Dosovitskiy et al., 2021) backbones of different scales, shown in Table 1. These backbones are chosen for several reasons: BERT and ResNet are commonly employed in prior works, which enables consistent comparisons. In addition, ViT and T5 are commonly used backbones for vision and language tasks, but relatively less experimented in prior debiasing works, which allows investigating the generalizability of our findings beyond existing research. Finally, each backbone is associated with a series of publicly available pre-trained checkpoints of different scales, with consistent network architecture and pre-training data, which enables cross-scale distillation and comparisons.

Table 1. Different scales of backbones in our experiments. h and d denote number of hidden layers and size of hidden dimension respectively. T, S, M, B, L refer to Tiny, Small, Medium, Base and Large version of the backbone. See Appendix for more details.

SCALE	В	ERT		Т5	RES	SNET	VIT		
Denilli	h	d	h	d	$\ h$	d	h	d	
Т	2	128	4	256	18	512	12	192	
S	4	256	8	384	34	512	12	384	
Μ	8	512	16	512	50	2048	12	768	
В	12	768	24	768	101	2048	24	1024	
L	24	1024	48	1024	152	2048	32	1280	



Figure 1. Framework for the analysis of distillability of debiasing methods. (a) training from scratch: we train a debiasing method M_i from scratch without knowledge distillation on different scales (teacher T and student S such that T > S) to obtain models f_T , f_S respectively. (b) Training with knowledge distillation: we apply knowledge distillation to transfer knowledge from teacher (f_T) to student $(g_{T->S})$. (c) Assessment: C1 determines if knowledge distillation in training a debiased model, and C3 compares different debiasing methods and backbones under knowledge distillation.

Evaluation To provide a comprehensive evaluation of robustness against spurious correlations, we compare the teacher $f_{\mathcal{T}}$ and the student $g_{\mathcal{T}->\mathcal{S}}$ from the following perspectives:

110 111

112

113

114

115

116

117

118 119

120

121

122

123

124

125 126

127

128

129

130

131

132

133

134

135

- **Performance on in-domain test set (ID,** ↑): This is the average performance on in-domain test set. A robust model should achieve high performance on this set to demonstrate general capability.
- Performance on out-of-domain test set / worst-136 137 group samples (OOD, \downarrow): For text datasets, we evalu-138 ate models on OOD test sets, comprised with specially 139 crafted hard samples (McCoy et al., 2019). Such sam-140 ples require real task-related signals to predict, where 141 biased models fall short. For image datasets, sam-142 ples are divided into groups based on their label and 143 spurious attributes. The worst performance on all sub-144 groups indicates the robustness to spurious correla-145 tions (Yang et al., 2023).
- Spurious gap (Spu. Gap, ↓): This metric calculates the performance gap between ID and OOD, which quantifies a model's vulnerability to spurious correlations. Ideally, a robust model should have high performances on both ID and OOD with a small spurious gap.
- 151 • Centralized Kernel Alignment (CKA) CKA is a 152 commonly adopted technique to measure the similarity 153 between activation matrices or hidden representations 154 of neural networks (Kornblith et al., 2019; Cortes 155 et al., 2012). Following previous work (Raghu et al., 156 2021; Nguyen et al., 2021), we use CKA by first 157 probing the intermediate representations from each 158 layer and then comparing all pairwise similarities 159 between representations of the teacher and student 160 models, under linear kernel CKA. 161

Similarly, we compare KD and Non-KD as above. We compute F1 score on QQP and accuracy on other datasets.

Datasets We use the following datasets

- CelebA (Liu et al., 2015) consists of 16k images of celebrity faces, where the objective is to predict "Blond_Hair" given "Male" as a spurious attribution.
- Waterbird (Sagawa et al., 2020) consists of synthetic images of birds from CUB dataset (Wah et al., 2011) and backgrounds (land & water) from Places (Zhou et al., 2018) dataset. The objective is to correctly infer "land bird" or "water bird," given the background as misleading information.
- MNLI (Williams et al., 2018) consists of 39k natural language inference (NLI) samples from various domains, where the objective is to classify relationship between a premise and a hypothesis as "Entailment", "Contradiction", or "Neutral". Previous studies discover that models are prone to negation words, lexical overlap, and sub-sequence biases in NLI task (Naik et al., 2018; Mendelson & Belinkov, 2021). We use HANS (McCoy et al., 2019) as the out-of-distribution test set (OOD) and SNLI (Bowman et al., 2015) as the transfer test set (Transfer), detailed below.
- **QQP** (Sharma et al., 2019) is a paraphrase identification (PI) dataset with 43k samples, where the objective is to predict if two questions are paraphrases of each other. Similar to MNLI, models are likely to be mislead by lexical overlap between two questions. We exploit PAWS (Zhang et al., 2019) as the out-of-distribution test set (OOD) and MRPC (Dolan & Brockett, 2005) as the transfer test set (Transfer), detailed below.

Debiasing Methods Experiments are conducted on a comprehensive list of commonly used debiasing methods, each of which is designed with special formulation and assumptions. We use (a) Empirical Risk Minimization (ERM) (standard training without debiasing techniques, (b) HypothesisOnly-PoE (Karimi Mahabadi et al.,

165																	
166	÷	ID						OOD				Spu. Gap					
167	Stu	$\mathcal{T}\!=\!T$	-0.5		ò		10	0.6			5	10	-1.1			5 o	5
107	÷.	$\mathcal{T}\!=\!S$	7.5	-0.2				1.5	0.0				6.0	-0.2			
168	ea	$\mathcal{T}\!=M$	9.9	2.1	-0.4			3.0	1.2	0.1			6.9	0.8	-0.5		
169	E	$\mathcal{T}\!=B$	12.5	4.3	1.9	-0.5		7.9	6.3	4.5	0.5		4.6	-2.0	-2.6	-1.0	
170	Gaj	$\mathcal{T} \!=\! L$	12.4	5.7	3.3	0.3	-0.6	13.3	12.9	11.0	6.4	0.1	-0.9	-7.3	-7.7	-6.1	-0.7
1/0			S = T	S = S	S = M	S = B	S = L	S = T	S = S	<i>S</i> = M	<i>S</i> = B	S = L	S = T	S = S	S = M	S = B	S = L
171																	
172																	

173Figure 2. C1: Teacher vs. Student: average performance gaps174between teacher and student models on ID, OOD, and Spurious175Gap across text datasets. X-axis and Y-axis show the scale of stu-176dent (S) and teacher (T) respectively: tiny (T), small (S), medium177(M), big (B), and large (L). Each cell shows the performance gap178between corresponding scales of a teacher and a student. See Appendix D for detailed results.

180 2020), (c) WeakLearner-PoE (Sanh et al., 2021), (d) Ker-181 nelWhitening (Gao et al., 2022), (e) AttentionPoE (Wang 182 et al., 2023), (f) σ-Damp (Puli et al., 2023), (g) Deep-183 FeatReweight (Kirichenko et al., 2023), and (h) PerSam-184 pleGrad (Ahn et al., 2023). The above debiasing methods 185 have a wide coverage of existing algorithms, ranging from 186 auxiliary biased model-based debiasing, to disentanglement 187 of representations. Meanwhile, they can handle multiple 188 types of shortcuts at the same time, without overfitting to 189 a specific bias. Details of these methods are provided in 190 Appendix A. 191

4. RQ1: Distillability of Debiasing Methods

193

195

196

We first examine if KD can effectively distill the debiasing capability from teachers to students of different scales.

197 Students become more biased than teachers We ob-198 serve that teachers consistently achieve better performance 199 than their smaller scale students on ID and OOD test sets 200 after knowledge distillation. The positive values in Fig-201 ure 2 show that although KD encourages students to mimic 202 their teachers in the logit space, it may undesirably increase student's susceptibility to spurious correlations in datasets 204 as the teacher's in-domain and debiasing capabilities are 205 not effectively transferred to the student. The prediction 206 agreement between teacher and student models show similar trend, where the student generally aligns with the teacher 208 on ID but often largely diverges from the teacher on OOD. 209 Furthermore, the extent of knowledge loss after distillation 210 varies depending on the relative scales of the teacher and 211 student models. For example, as depicted in Figure 2, when 212 S is tiny (S = T), more debiasing power is lost, shown by 213 mostly positive values in spurious gap. When \mathcal{T} is large (\mathcal{T} 214 = L), more ID knowledge is lost, shown by mostly negative 215 values in spurious gap. The results show that if a teacher 216 model learns a partially debiased representation but still 217 retains residual biases, the student might amplify this bias 218 rather than mitigate it. 219



Figure 3. **C1: Teacher vs. Student**: density of predicted probability on text datasets. On OOD, students has larger deviation in prediction confidence than teachers. See Appendix D for detailed results.

Students show diverse distribution shifts in predicted probabilities To understand the influence of KD on the debiasing capabilities of students, we investigate the output probability distribution $P_{\mathcal{C}}(y=1)$. Our findings show that KD significantly alter the predicted probability distribution, despite its training objective of matching output logits. This perturbation is often larger on OOD than ID test sets, which explains the larger performance drop observed in students compared to their teachers on OOD, as illustrated in Figure 3. We also observed that teachers tend to provide slightly more confident predictions on ID while more moderate predictions on OOD. Such behavior is not successfully transferred to students through KD. Such distinct behaviors on different samples may encourage models to overfit to data distributions of the training sets or to over-represented groups, which can effectively amplify reliance on shortcuts over robust features. In addition, the training sets of teachers often contain biased examples or do not equally represent all sub-groups, which leads students to inherit and potentially amplify these biases. Consequently, students often perform worse than their teachers on OOD.

Students and teachers show different attention on ID and OOD To have a deeper understanding of the teacherstudent divergence, we further probe the internal representations when making predictions on ID and OOD data. Results shows that after distillation, students try to mimic teachers on ID (left). The earlier layers of students follow earlier layers of teachers, and similarly mid and later layers. This indicates that KD can transfer knowledge of ID data from the larger teachers into smaller students. Pn OOD (right), however, we observe similar pattern but it is not fully preserved. In particular, it is challenging for the mid and later layers of the students to follow closely to those of the teachers, which explains the performance degradation on OOD after KD, see Figure 5.

Potential for new debiasing capabilities for students beyond teacher abilities We compare prediction agreement between teacher and student models. When T is large (T



Figure 4. C1: Teacher vs. Student: prediction agreement on text datasets. Left: varying S (X-axis) given a fixed large teacher (T = L). Right: varying T (X-axis) given a fixed tiny student (S = T). Agreement increases as the scale of teacher and student get closer. See Appendix D for detailed results.

= L), we observe an increase in prediction agreement as Sscales up, with consistently higher agreement on ID than OOD, as shown in the left plot in Figure 4. Conversely, when S is tiny (S = T), the prediction agreement diminishes as \mathcal{T} scales up, with higher agreement on OOD than ID, the right plot in Figure 4. The imperfect agreement between teacher and student contradicts with the foundational assumptions of knowledge distillation, which assumes that students should closely mimic their teachers. However, interestingly, this unexpected behavior may not always lead to performance degradation. Sometimes it enables students to generalize to out of domain data. In particular, there are instances where students make correct predictions where their teachers do not, see the left plot in Figure 10. Students can sometimes outperform their teachers perhaps because they may learn additional patterns during the knowledge distillation process, which allows them to generalize better than their teachers. The above result suggests that students may sometimes acquire debiasing capabilities that surpass those of their teachers, which we believe is a novel avenue for robust model training.

Larger teachers do not guarantee more robust students Our findings show that a more capable teacher does not guarantee a less biased student in debiasing tasks. With a fixed student scale (as seen in the columns of Figure 2), increasing the teacher's scale does not consistently reduce performance gap or spurious gap. Sometimes, a larger teacher may degrade the debiasing capability of the student. For example, when S = T, increasing the teacher scale from M to B increases the spurious gap from 6.5 to 8.1 on ERM, i.e. a more biased model. Moreover, when when S = T, increasing Tresult in a drop of teacher-student agreement, indicating that the student fails to follow the teacher, see right plot in Figure 4. We attribute this finding to two reasons. Firstly, the capabilities of students are substantially bounded by their scale, and using a much larger and capable teacher may exceed the student's capacity for effective learning (Cho & Hariharan, 2019). Secondly, training students with debias-



Figure 5. **C1: Teacher vs. Student**: Centered Kernel Alignment on ID (left) and OOD (right). Highers values indicate higher similarity. X-axis and Y-axis refer to the layers of teacher (T) and student (S) respectively. See Appendix D for detailed results.

ing objectives and knowledge distillation at the same time results in optimization problem, which may trap students' parameters in local optima and affect their robustness to spurious correlations.

Students with similar scales to their teachers learn better The effectiveness of debiasing ability transfer through distillation is greatly affected by the scale similarity between teacher and student. As the teacher and the student become similar in scale (near the diagonal cells in Figure 2), the differences on test set performance and spurious gaps decrease. However, a larger mismatch in scale (far from diagonal) results in more pronounced differences, see Figure 2. Similarly, the teacher-student agreement increases as T and Salign more closely, see Figure 4. This is likely because models of similar scales have comparable expressive power and extracts similar features, which can lead to more effective knowledge transfer, better bias mitigation, and higher prediction agreement.

5. RQ2: Distillation vs. Standard Training

We asses if training with knowledge distillation (KD) can improve a model's debiasing performance compared to standard training (Non-KD).

Non-KD is less biased than KD Our results show that debiasing models trained from scratch (Non-KD) have lower ID performance than those trained with KD. However, the Non-KD models achieve almost no changes on OOD, leading to smaller spurious gaps, see Figure 6. We hypothesize that the distillation objective of matching logits, despite effective on ID, may potentially inject additional spurious correlations and distract the model from prioritizing robust features, as the teacher is not fully unbiased.

KD does not improve generalization An interesting finding is that both Non-KD and KD have similar average prediction agreements on both ID and OOD. However,



Figure 6. C2: Non-KD vs. KD: average performance gap between teacher and student on ID, OOD, and Spurious Gap across text datasets. X-axis and Y-axis show the scale of student (S) and teacher (T) respectively: : tiny (T), small (S), medium (M), big (B), and large (L). Each cell shows the performance gap of the corresponding scales of a teacher and a student. See Appendix D for detailed results.

283

284

285

286

287

288

289

299

300

301

302

319

320

322

323 324

325

327

328

329



Figure 7. **C2: Non-KD vs. KD**: density of predicted probability on text datasets. On OOD, KD has larger deviation in prediction confidence than Non-KD. See Appendix D for detailed results.

303 the agreement on OOD varies significantly depending on 304 dataset, debiasing method, and backbone model. This sug-305 gests that training solely with the original data (Non-KD) 306 is sufficiently effective for debiasing, and introducing exter-307 nal knowledge via KD does not yield significant improve-308 ments. This result can be attributed to KD's impact on 309 model confidence; models trained with KD tend to produce more confident predictions than models trained without KD, 311 see Figure 7, which is key to degenerate performance on 312 OOD (Utama et al., 2020; Sanh et al., 2021). Such overcon-313 fidence could be a critical factor in degraded performance 314 on OOD tasks. Moreover, such minimal contribution of KD 315 remains unchanged even when stronger external knowledge 316 (a larger teacher) or a more capable learner (a larger student) 317 is used, see Figure 8. 318

6. RQ3: Effect of Debiasing Methods

We assess the effect of different debiasing methods and backbones on our earlier findings.

Different transfer patterns across methods Our results show that the transfer patterns are heavily influenced by the formulation of debiasing method. For example, logitbased PoE methods, such as HypothesisOnly-PoE and WeakLearner-PoE, show similar trends in performance



Figure 8. **C2:** Non-KD vs. KD: prediction agreement on text datasets. Left: varying S (X-axis) given a fixed teacher with T = L. Right: varying T (X-axis) given a fixed student with T = T. Agreement does not increase significantly as T and S scale up. See Appendix D for detailed results.

changes and spurious gaps, in contrast to the representation disentanglement method (KernelWhitening), see Figure 9.

Sensitivity to backbones The distillability of KD appears to varies with the architecture of the backbones and randomness in the training. KernelWhitening and WeakLearner-PoE are two methods particularly sensitive to the scale of backbone and random seeds, which controls factors such as data sampling and ordering.

Robustness to different biases transfer differently We observe that OOD and Transfer show different transfer patterns, where performance gap on Transfer exhibit much larger variations the student scales up, see Figure 9. This suggests that smaller students may outperform larger ones on OOD, indicating that during KD, larger students may become more prone to certain biases (OOD) but more resilient to others.

Universal transfer patterns in debiasing methods A number of debiasing methods show consistent changes in robustness after KD, which suggest the potential for an empirical universal transfer pattern. Specifically on text datasets, the performance gap between teacher and student models on OOD and Transfer Spurious Gap fall in the range of [0, 5] and [-5, 0] respectively, see Figure 9. Such change in performance is consistent across different scales of \mathcal{T} and \mathcal{S} , which allows for predictable performance after KD. Similarly, the performance gap between models trained using Non-KD and KD remains stable on OOD, falling in range [-1, 1] across different scales.

7. Potential Solutions

Based on the above analyses, we summarize the key findings on distilling debiasing models as follows:

• Training distribution significantly affect successful distillation of debiasing capabilities.

372

- Student models with similar scale to their teachers can better obtain debiasing knowledge from their teachers.
- The objectives of KD may introduce additional optimization challenges, especially with the presence of debiasing objectives.

To further improve the distillability of debiasing methods, we propose three solutions:

Data augmentation (DA) There is broad evidence that models becomes biased by relying on spurious features in the training set (Wu et al., 2022; Ahn et al., 2023), which is amplified by misrepresentation of specific classes or labeling errors. Prior studies have highlighted the important role of data in knowledge distillation (Stanton et al., 2021). Based on these prior studies and our findings, we hypothesize that providing high quality data and augmenting data size can improve the process of distilling the debiasing capability from teacher to student. For text datasets, we employ the data generated by Seq-Z filtering (Wu et al., 2022) as training set for both teacher and student models. For image datasets, we employ training and validation sets where the sub-groups are equally represented (Kirichenko et al., 2023).

Iterative knowledge distillation (IKD) Our results indicate that the transfer of debiasing capabilities is more effective between teachers and students of similar scales. Therefore, we propose to leverage Iterative Knowledge Distillation (IKD) (Liu et al., 2023): given a teacher of scale 360 361 S_N , we first distill it to a student of scale S_{N-1} , where \mathcal{S}_{N-1} is the closest neighbor of \mathcal{S}_N in scale. Then the 362 newly-distilled student acts as a teacher and transfer the 363 knowledge to a model of smaller scale S_{N-2} , where S_{N-2} is the closest neighbor of S_{N-1} in scale. We repeat this 365 process iteratively by gradually decreasing student scale, 367 such that the knowledge can be transferred smoothly from 368 a large scale model to a small scale model. This step-wise 369 approach enables a smooth knowledge transfer from larger 370 to smaller scale models, and potentially improves debiasing effectiveness at each step. 371

373 Initialize student with teacher weights (Init) Previous 374 research by Stanton et al. (2021) has discovered that initial-375 izing a student model with the weights of its teacher can 376 increase their centered kernel alignment (Kornblith et al., 377 2019) in activation space. This approach can head-start the 378 student model with a stronger debiasing capability from the 379 teacher. It can also help alleviate potential optimization ob-380 stacles and stuck in local optima. If the teacher and student 381 models are of the same scale, we initialize the student with 382 the teacher parameters. If the teacher is larger, we initialize 383 the student with the first few layers of the teacher. 384

Table 2. Improvement of distillability. Vanilla refers to standard knowledge distillation, DA, IKD, and Init represent data augmentation, iterative knowledge distillation, and initialization of student with teacher weights (Init) as our three solutions to improve the distillability of debiasing methods.

DIFF IN	ID (↓)	OOD (↓)	Spu. Gap (↓)	ID (↓)	OOD (↓)	Spu. Gap (↓)
	TEAG	CHER - S	Student	∥ Nor	• KD - 1	KD
VANILLA	5.1	7.3	12.7	1.4	0.7	2.2
+ DA + IKD	2.3	5.4 5.9	8.2 9.9	0.2	0.2 0.5	0.5 1.6
+ Init.	4.7	6.5	11.5	1.3	0.7	2.0

Results Table 2 shows that all three solutions result in improved distillability. Specifically, on spurious gap between teacher and student, data augmentation (DA), iterative knowledge distillation (IKD), and initialization with teacher weights (Init) yield performance gains of 4.5, 2.8, 1.2 absolute points compared to vanilla KD across datasets and backbones respectively. On spurious gap between Non-KD and KD, DA, IKD, and Init outperforms vanilla KD by 1.7, 0.6 and 0.2 absolute points respectively. We find that DA achieves the largest improvement, since the root cause of spurious correlations come from the underlying dataset (Chen et al., 2018). Debiasing the dataset itself can benefit all training methods including knowledge distillation. As noted by previous work (Stanton et al., 2021), Init may facilitate teacher-student agreement in activation space, but result in non-significant gains, which aligns with our findings as well.

8. Related Work

Bias mitigation in NLU: Debiasing approaches usually employ a biased model to inform the training of a robust model (Clark et al., 2019; Karimi Mahabadi et al., 2020; Sanh et al., 2021; Utama et al., 2020; Cheng et al., 2024). Other methods aim at learning debiased or robust representations (Gao et al., 2022; Lyu et al., 2022; Wang et al., 2023; Jeon et al., 2023; Reif & Schwartz, 2023), or removing bias-encoding parameters (Meissner et al., 2022; Yu et al., 2023). Other works include measurement of bias of specific words with statistical test (Gardner et al., 2021), generating non-biased samples (Wu et al., 2022), identification of biasencoding parameters (Yu et al., 2023), when bias mitigation works (Ravichander et al., 2023), and bias transfer from other models (Jin et al., 2021).

Bias mitigation in vision: In vision, worst-group performance is measured as a sign of model robustness. Several works investigates how to learn debiased models from failure cases (Nam et al., 2020), biased representations (Bahng



Figure 9. **C3: Comparison Between Debiasing Methods**: performance gap between Teacher and Student (Above), Non-KD and KD (Lower) on text datasets. Detailed results are shown in Appendix.



399

400

408

409

410

411

Figure 10. **C1: Teacher vs. Student** (Left) and **C2: Non-KD vs. KD** (Right): correctly predicted examples on OOD on text datasets.

et al., 2020), multiple biased models (Kim et al., 2022), and 412 by simply re-training the last layer of a neural model (i.e. 413 414 the classification layer) with additional equally represented data (Kirichenko et al., 2023; LaBonte et al., 2023). Li et al. 415 (2023) showed that multiple spurious features can occur 416 in a dataset, while suppressing one may inevitably boost 417 another one. Other perspectives for debiasing include causal 418 attention (Wang et al., 2021), building uniform margin clas-419 420 sifiers (Puli et al., 2023), using representations from earlier layers (Tiwari et al., 2024), and neural collapse (Wang et al., 421 422 2024), where feature space collapses into a stable geometric 423 structure that results in robustness and generalizability.

424 Knowledge distillation: Knowledge Distillation (KD) is 425 initially proposed to transfer knowledge from a larger model 426 (teacher) to a smaller model (student), by encouraging the 427 student to follow the teacher on prediction logits (Hinton 428 et al., 2015), learned features (Romero et al., 2015; Wang 429 et al., 2020), attention map (Zagoruyko & Komodakis, 2017; 430 Chen et al., 2021), activation patterns (Huang & Wang, 431 2017; Heo et al., 2019). Later works discovered that KD can 432 be viewed as a special form of regularization similar to label 433 smoothing (Szegedy et al., 2016), providing no task-specific 434 knowledge. However, on text classification tasks, whether 435 KD can regularize the student depends on the choice of 436 teacher model (Sultan, 2023), which may result in opposite 437 model confidence between teacher and student compared to 438 label smoothing. Stanton et al. (2021) discovers that opti-439

mization and dataset details are crucial to matching students to teachers, and such matching does not guarantee better generalization ability of students. Xue et al. (2023) investigates cross-modal KD, where the teacher functions on a different modality or extra modalities than student. The authors proposes modality fusing hypothesis, which claims that modality decisive features are critical for the effectiveness of cross-modal KD. However, despite briefly discussed (Cho & Hariharan, 2019; Tiwari et al., 2024), the potential of knowledge distillation to transfer debiasing capabilities across different modalities and backbone models remains underexplored and poorly understood in existing work.

9. Conclusion

We present the first study on the distillability of debiasing capabilities between neural models, and the extent of bias transfer through knowledge distillation (KD). We evaluate eight popular debiasing methods and five scales of backbones on four datasets. Extensive experiments show that vanilla KD does not consistently preserve debiasing capabilities; in many cases, student models become more reliant on spurious correlations than their teachers; the effectiveness of debiasing transfer depends on model scale similarity-distillation works best when teacher and student models are comparable in complexity; and larger teachers do not always yield more robust students, which indicates the need for targeted debiasing strategies in KD. We propose three solutions-data augmentation, iterative KD, and student initialization-which significantly improve the distillability of debiasing methods and contribution of KD on debiasing.

In future we will investigate self-distilled debiasing, where the student iteratively distills knowledge from itself rather than relying on a fixed teacher. A potential improvement is to explicitly guide the student using counterfactual data augmentation during distillation.

440 Impact Statement

441

442

443

444

445

446

447

448

449

450

Our research focuses on mitigating dataset biases in text and vision datasets, and understanding why debiasing methods may fail under knowledge distillation. The broader impacts of our work are in advancing dataset fairness and potentially enhancing decision-making based on data. Our work contributes to improving the accuracy and reliability of NLP and vision models, as well as their trust and adoption.

References

- Ahn, S., Kim, S., and Yun, S.-Y. Mitigating dataset bias by using per-sample gradient. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum? id=7mgUec-7GMv.
- 457 Bahng, H., Chun, S., Yun, S., Choo, J., and Oh, S. J. 458 Learning de-biased representations with biased repre-459 sentations. In III, H. D. and Singh, A. (eds.), Pro-460 ceedings of the 37th International Conference on Ma-461 chine Learning, volume 119 of Proceedings of Machine 462 Learning Research, pp. 528-539. PMLR, 13-18 Jul 463 2020. URL https://proceedings.mlr.press/ 464 v119/bahng20a.html. 465
- Bowman, S. R., Angeli, G., Potts, C., and Manning, C. D. 466 A large annotated corpus for learning natural language 467 inference. In Màrquez, L., Callison-Burch, C., and Su, J. 468 (eds.), Proceedings of the 2015 Conference on Empirical 469 Methods in Natural Language Processing, pp. 632–642, 470 Lisbon, Portugal, September 2015. Association for Com-471 putational Linguistics. doi: 10.18653/v1/D15-1075. URL 472 https://aclanthology.org/D15-1075/. 473
- Chen, D., Mei, J.-P., Zhang, Y., Wang, C., Wang, Z., Feng,
 Y., and Chen, C. Cross-layer distillation with semantic calibration. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(8):7028–7036, May 2021. doi: 10.
 1609/aaai.v35i8.16865. URL https://ojs.aaai.
 org/index.php/AAAI/article/view/16865.
- Chen, I., Johansson, F. D., and Sontag, D. Why is my
 classifier discriminatory? *Advances in neural information processing systems*, 31, 2018.
- 484 Cheng, J. and Amiri, H. FairFlow: Mitigating dataset 485 biases through undecided learning for natural lan-486 guage understanding. In Al-Onaizan, Y., Bansal, 487 M., and Chen, Y.-N. (eds.), Proceedings of the 2024 488 Conference on Empirical Methods in Natural Lan-489 guage Processing, pp. 21960–21975, Miami, Florida, 490 USA, November 2024. Association for Computa-491 tional Linguistics. doi: 10.18653/v1/2024.emnlp-main. 492 1225. URL https://aclanthology.org/2024. 493 emnlp-main.1225/. 494

- Cheng, J., Elgaar, M., Vakil, N., and Amiri, H. Cognivoice: Multimodal and multilingual fusion networks for mild cognitive impairment assessment from spontaneous speech. *arXiv preprint arXiv:2407.13660*, 2024.
- Chew, O., Lin, H.-T., Chang, K.-W., and Huang, K.-H. Understanding and mitigating spurious correlations in text classification with neighborhood analysis. In Graham, Y. and Purver, M. (eds.), *Findings of the Association for Computational Linguistics: EACL 2024*, pp. 1013–1025, St. Julian's, Malta, March 2024. Association for Computational Linguistics. URL https://aclanthology. org/2024.findings-eacl.68/.
- Cho, J. H. and Hariharan, B. On the efficacy of knowledge distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.
- Clark, C., Yatskar, M., and Zettlemoyer, L. Don't take the easy way out: Ensemble based methods for avoiding known dataset biases. In Inui, K., Jiang, J., Ng, V., and Wan, X. (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 4069–4082, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/ D19-1418. URL https://aclanthology.org/ D19-1418/.
- Cortes, C., Mohri, M., and Rostamizadeh, A. Algorithms for learning kernels based on centered alignment. *Journal of Machine Learning Research*, 13(28): 795–828, 2012. URL http://jmlr.org/papers/v13/cortes12a.html.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In Burstein, J., Doran, C., and Solorio, T. (eds.), Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/N19-1423/.
- Dolan, W. B. and Brockett, C. Automatically constructing a corpus of sentential paraphrases. In *Proceedings* of the Third International Workshop on Paraphrasing (*IWP2005*), 2005. URL https://aclanthology. org/105-5002/.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby,

N. An image is worth 16x16 words: Transformers for
image recognition at scale. In *International Conference on Learning Representations*, 2021. URL https://
openreview.net/forum?id=YicbFdNTTy.

- 499 Gao, S., Dou, S., Zhang, Q., and Huang, X. Kernel-500 whitening: Overcome dataset bias with isotropic sen-501 tence embedding. In Goldberg, Y., Kozareva, Z., 502 and Zhang, Y. (eds.), Proceedings of the 2022 Con-503 ference on Empirical Methods in Natural Language 504 Processing, pp. 4112–4122, Abu Dhabi, United Arab 505 Emirates, December 2022. Association for Computa-506 tional Linguistics. doi: 10.18653/v1/2022.emnlp-main. 507 275. URL https://aclanthology.org/2022. 508 emnlp-main.275/. 509
- 510 Gardner, M., Merrill, W., Dodge, J., Peters, M., Ross, 511 A., Singh, S., and Smith, N. A. Competency prob-512 lems: On finding and removing artifacts in language 513 data. In Moens, M.-F., Huang, X., Specia, L., and 514 Yih, S. W.-t. (eds.), Proceedings of the 2021 Confer-515 ence on Empirical Methods in Natural Language Pro-516 cessing, pp. 1801–1813, Online and Punta Cana, Domini-517 can Republic, November 2021. Association for Computa-518 tional Linguistics. doi: 10.18653/v1/2021.emnlp-main. 519 135. URL https://aclanthology.org/2021. 520 emnlp-main.135/. 521
- Guo, Q., Tang, Y., Ouyang, Y., Wu, Z., and Dai, 522 Debias NLU datasets via training-free pertur-X. 523 In Bouamor, H., Pino, J., and Bali, K. 524 bations. (eds.), Findings of the Association for Computational 525 Linguistics: EMNLP 2023, pp. 10886-10901, Singa-526 pore, December 2023. Association for Computational 527 Linguistics. doi: 10.18653/v1/2023.findings-emnlp. 528 726. URL https://aclanthology.org/2023. 529 530 findings-emnlp.726/.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (*CVPR*), June 2016.
- 536 Heo, B., Lee, M., Yun, S., and Choi, J. Y. Knowl-537 edge transfer via distillation of activation boundaries 538 formed by hidden neurons. Proceedings of the 539 AAAI Conference on Artificial Intelligence, 33(01): 540 3779-3787, Jul. 2019. doi: 10.1609/aaai.v33i01. 541 33013779. URL https://ojs.aaai.org/index. 542 php/AAAI/article/view/4264. 543
- Hinton, G., Vinyals, O., and Dean, J. Distilling
 the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Hinton, G. E. Training products of experts by minimizing contrastive divergence. *Neural Comput.*, 14(8):

1771-1800, aug 2002. ISSN 0899-7667. doi: 10.1162/ 089976602760128018. URL https://doi.org/10. 1162/089976602760128018.

- Huang, Z. and Wang, N. Like what you like: Knowledge distill via neuron selectivity transfer. *arXiv preprint arXiv:1707.01219*, 2017.
- Jeon, E., Lee, M., Park, J., Kim, Y., Mok, W.-L., and Lee, S. Improving bias mitigation through bias experts in natural language understanding. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 11053–11066, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023. emnlp-main.681. URL https://aclanthology. org/2023.emnlp-main.681/.
- Jin, X., Barbieri, F., Kennedy, B., Mostafazadeh Davani, A., Neves, L., and Ren, X. On transferability of bias mitigation effects in language model fine-tuning. In Toutanova, K., Rumshisky, A., Zettlemoyer, L., Hakkani-Tur, D., Beltagy, I., Bethard, S., Cotterell, R., Chakraborty, T., and Zhou, Y. (eds.), Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 3770–3783, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021. naacl-main.296. URL https://aclanthology. org/2021.naacl-main.296/.
- Karimi Mahabadi, R., Belinkov, Y., and Henderson, J. Endto-end bias mitigation by modelling biases in corpora. In Jurafsky, D., Chai, J., Schluter, N., and Tetreault, J. (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 8706–8716, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.769. URL https: //aclanthology.org/2020.acl-main.769/.
- Kim, N., HWANG, S., Ahn, S., Park, J., and Kwak, S. Learning debiased classifier with biased committee. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), Advances in Neural Information Processing Systems, volume 35, pp. 18403–18415. Curran Associates, Inc., 2022. URL https://proceedings.neurips. cc/paper_files/paper/2022/file/ 750046157471c56235a781f2eff6e226-Paper-Conference pdf.
- Kirichenko, P., Izmailov, P., and Wilson, A. G. Last layer re-training is sufficient for robustness to spurious correlations. In *The Eleventh International Conference on Learning Representations*, 2023. URL https: //openreview.net/forum?id=Zb6c8A-Fghk.

Kornblith, S., Norouzi, M., Lee, H., and Hinton, G. Similarity of neural network representations revisited. In
Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 3519–3529. PMLR, 09–15 Jun 2019.
URL https://proceedings.mlr.press/v97/

557 kornblith19a.html.

- 558 LaBonte, T., Muthukumar, V., and Kumar, A. Towards 559 last-layer retraining for group robustness with fewer 560 annotations. In Oh, A., Naumann, T., Globerson, 561 A., Saenko, K., Hardt, M., and Levine, S. (eds.), 562 Advances in Neural Information Processing Systems, 563 volume 36, pp. 11552-11579. Curran Associates, Inc., 564 2023. URL https://proceedings.neurips. 565 cc/paper files/paper/2023/file/ 566 265bee74aee86df77e8e36d25e786ab5-Paper-Conference.116/. 567
- 568

pdf.

Li, Z., Evtimov, I., Gordo, A., Hazirbas, C., Hassner, T.,
Ferrer, C. C., Xu, C., and Ibrahim, M. A whac-a-mole dilemma: Shortcuts come in multiples where mitigating one amplifies others. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (*CVPR*), pp. 20071–20082, June 2023.

Liu, J., Wang, P., Shang, Z., and Wu, C. Iterde: An iterative knowledge distillation framework for knowledge graph embeddings. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37:4488–4496, Jun. 2023. doi: 10.
1609/aaai.v37i4.25570. URL https://ojs.aaai.org/index.php/AAAI/article/view/25570.

- Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning
 face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pp. 3730–3738, 2015.
- Lyu, Y., Li, P., Yang, Y., de Rijke, M., Ren, P., Zhao, Y., Yin,
 D., and Ren, Z. Feature-level debiased natural language understanding. In *Proceedings of the AAAI Conference* on Artificial Intelligence, 2022.
- McCoy, T., Pavlick, E., and Linzen, T. Right for the 592 wrong reasons: Diagnosing syntactic heuristics in nat-593 ural language inference. In Korhonen, A., Traum, D., 594 and Màrquez, L. (eds.), Proceedings of the 57th Annual 595 Meeting of the Association for Computational Linguis-596 tics, pp. 3428-3448, Florence, Italy, July 2019. Associ-597 ation for Computational Linguistics. doi: 10.18653/v1/ 598 P19-1334. URL https://aclanthology.org/ 599 P19-1334/. 600
- Meissner, J. M., Sugawara, S., and Aizawa, A. Debiasing masks: A new framework for shortcut mitigation in NLU. In Goldberg, Y., Kozareva, Z.,

and Zhang, Y. (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 7607–7613, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main. 517. URL https://aclanthology.org/2022. emnlp-main.517/.

- Mendelson, M. and Belinkov, Y. Debiasing methods in natural language understanding make bias more accessible. In Moens, M.-F., Huang, X., Specia, L., and Yih, S. W.-t. (eds.), Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pp. 1545–1557, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main. 116. URL https://aclanthology.org/2021. onfemtpomain.116/.
- Naik, A., Ravichander, A., Sadeh, N., Rose, C., and Neubig, G. Stress test evaluation for natural language inference. In Bender, E. M., Derczynski, L., and Isabelle, P. (eds.), Proceedings of the 27th International Conference on Computational Linguistics, pp. 2340–2353, Santa Fe, New Mexico, USA, August 2018. Association for Computational Linguistics. URL https://aclanthology. org/C18-1198/.
- Nam, J., Cha, H., Ahn, S., Lee, J., and Shin, J. Learning from failure: De-biasing classifier from biased classifier. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), Advances in Neural Information Processing Systems, volume 33, pp. 20673–20684. Curran Associates, Inc., 2020. URL https://proceedings.neurips. cc/paper_files/paper/2020/file/ eddc3427c5d77843c2253fle799fe933-Paper. pdf.
- Nguyen, T., Raghu, M., and Kornblith, S. Do wide and deep networks learn the same things? uncovering how neural network representations vary with width and depth. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum? id=KJNcAkY8tY4.
- Puli, A. M., Zhang, L., Wald, Y., and Ranganath, R. Don't blame dataset shift! shortcut learning due to gradients and cross entropy. In Oh, A., Naumann, T., Globerson, A., Saenko, K., Hardt, M., and Levine, S. (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 71874–71910. Curran Associates, Inc., 2023. URL https://proceedings.neurips. cc/paper_files/paper/2023/file/ e35460304fdf6df523f068a59aaf8829-Paper-Conference pdf.

Raghu, M., Unterthiner, T., Kornblith, S., Zhang, C., and 605 606 Dosovitskiy, A. Do vision transformers see like convolu-607 tional neural networks? In Ranzato, M., Beygelzimer, 608 A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), 609 Advances in Neural Information Processing Systems, volume 34, pp. 12116–12128. Curran Associates, Inc., 610 611 2021. URL https://proceedings.neurips. 612 cc/paper_files/paper/2021/file/ 613 652cf38361a209088302ba2b8b7f51e0-Paper. 614 pdf. 615 Ravichander, A., Stacey, J., and Rei, M. When and why 616

does bias mitigation work? In Bouamor, H., Pino, J., and
Bali, K. (eds.), *Findings of the Association for Computa- tional Linguistics: EMNLP 2023*, pp. 9233–9247, Singapore, December 2023. Association for Computational
Linguistics. doi: 10.18653/v1/2023.findings-emnlp.
619. URL https://aclanthology.org/2023.
findings-emnlp.619/.

624 Reif, Y. and Schwartz, R. Fighting bias with bias: 625 Promoting model robustness by amplifying dataset bi-626 ases. In Rogers, A., Boyd-Graber, J., and Okazaki, 627 N. (eds.), Findings of the Association for Compu-628 tational Linguistics: ACL 2023, pp. 13169-13189, 629 Toronto, Canada, July 2023. Association for Computa-630 tional Linguistics. doi: 10.18653/v1/2023.findings-acl. 631 833. URL https://aclanthology.org/2023. 632 findings-acl.833/. 633

Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta,
C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *International Conference on Learning Representations*,
2015.

634

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

- Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P. Distributionally robust neural networks. In International Conference on Learning Representations, 2020. URL https://openreview.net/forum? id=ryxGuJrFvS.
- Sanh, V., Wolf, T., Belinkov, Y., and Rush, A. M. Learning from others' mistakes: Avoiding dataset biases without modeling them. In *International Conference on Learning Representations*, 2021. URL https://openreview. net/forum?id=Hf3qXoiNkR.
- Sharma, L., Graesser, L., Nangia, N., and Evci, U. Natural language understanding with the quora question pairs dataset. *arXiv e-prints*, 2019. URL https://arxiv. org/abs/1907.01041.
- Stanton, S. D., Izmailov, P., Kirichenko, P., Alemi, A. A., and Wilson, A. G. Does knowledge distillation really work? In Beygelzimer, A., Dauphin, Y., Liang,

P., and Vaughan, J. W. (eds.), Advances in Neural Information Processing Systems, 2021. URL https: //openreview.net/forum?id=7J-fKoXiReA.

- Sultan, M. Knowledge distillation ≈ label smoothing: Fact or fallacy? In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Meth*ods in Natural Language Processing, pp. 4469–4477, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main. 271. URL https://aclanthology.org/2023. emnlp-main.271/.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- Tay, Y., Dehghani, M., Rao, J., Fedus, W., Abnar, S., Chung, H. W., Narang, S., Yogatama, D., Vaswani, A., and Metzler, D. Scale efficiently: Insights from pretraining and finetuning transformers. In *International Conference on Learning Representations*, 2022. URL https: //openreview.net/forum?id=f20YVDyfIB.
- Tiwari, R., Sivasubramanian, D., Mekala, A., Ramakrishnan, G., and Shenoy, P. Using early readouts to mediate featural bias in distillation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* (*WACV*), pp. 2638–2647, January 2024.
- Turc, I., Chang, M.-W., Lee, K., and Toutanova, K. Wellread students learn better: On the importance of pretraining compact models, 2019.
- Utama, P. A., Moosavi, N. S., and Gurevych, I. Towards debiasing NLU models from unknown biases. In Webber, B., Cohn, T., He, Y., and Liu, Y. (eds.), Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 7597– 7610, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main. 613. URL https://aclanthology.org/2020. emnlp-main.613/.
- Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. Caltech-ucsd birds-200-2011 (cub-200-2011). Technical report, California Institute of Technology, 2011.
- Wang, F., Huang, J. Y., Yan, T., Zhou, W., and Chen, M. Robust natural language understanding with residual attention debiasing. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 504–519, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl. 32. URL https://aclanthology.org/2023. findings-acl.32/.

- Wang, T., Zhou, C., Sun, Q., and Zhang, H. Causal attention
 for unbiased visual recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 3091–3100, October 2021.
- Wang, X., Fu, T., Liao, S., Wang, S., Lei, Z., and Mei,
 T. Exclusivity-consistency regularized knowledge distillation for face recognition. In Vedaldi, A., Bischof,
 H., Brox, T., and Frahm, J.-M. (eds.), *Computer Vision – ECCV 2020*, pp. 325–342, Cham, 2020. Springer International Publishing. ISBN 978-3-030-58586-0.
- Wang, Y., Sun, J., Wang, C., Zhang, M., and Yang, M. Navigate beyond shortcuts: Debiased learning through the lens of neural collapse. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (*CVPR*), pp. 12322–12331, June 2024.
- Williams, A., Nangia, N., and Bowman, S. A broad-677 678 coverage challenge corpus for sentence understanding through inference. In Walker, M., Ji, H., and Stent, 679 A. (eds.), Proceedings of the 2018 Conference of the 680 North American Chapter of the Association for Com-681 putational Linguistics: Human Language Technologies, 682 Volume 1 (Long Papers), pp. 1112–1122, New Orleans, 683 Louisiana, June 2018. Association for Computational 684 Linguistics. doi: 10.18653/v1/N18-1101. URL https: 685 //aclanthology.org/N18-1101/. 686
- 687 Wu, Y., Gardner, M., Stenetorp, P., and Dasigi, P. Generat-688 ing data to mitigate spurious correlations in natural lan-689 guage inference datasets. In Muresan, S., Nakov, P., and 690 Villavicencio, A. (eds.), Proceedings of the 60th Annual 691 Meeting of the Association for Computational Linguistics 692 (Volume 1: Long Papers), pp. 2660-2676, Dublin, Ire-693 land, May 2022. Association for Computational Linguis-694 tics. doi: 10.18653/v1/2022.acl-long.190. URL https: 695 //aclanthology.org/2022.acl-long.190/. 696
- Ku, Z., Jin, R., Shen, B., and Zhu, S. Nystrom approximation for sparse kernel methods: Theoretical analysis and empirical evaluation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 29(1), Feb. 2015. doi: 10.1609/aaai.v29i1.9626. URL https://ojs.aaai.org/index.php/AAAI/article/view/9626.
- Xue, Z., Gao, Z., Ren, S., and Zhao, H. The modality focusing hypothesis: Towards understanding crossmodal knowledge distillation. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum? id=w0QXrZ3N-s.
- Yang, Y., Zhang, H., Katabi, D., and Ghassemi, M.
 Change is hard: A closer look at subpopulation shift.
 In Krause, A., Brunskill, E., Cho, K., Engelhardt,
 B., Sabato, S., and Scarlett, J. (eds.), *Proceedings of*

the 40th International Conference on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pp. 39584–39622. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/ v202/yang23s.html.

- Yu, C., Jeoung, S., Kasi, A., Yu, P., and Ji, H. Unlearning bias in language models by partitioning gradients. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 6032–6048, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl. 375. URL https://aclanthology.org/2023. findings-acl.375/.
- Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum? id=Sks9_ajex.
- Zhang, Y., Baldridge, J., and He, L. PAWS: Paraphrase adversaries from word scrambling. In Burstein, J., Doran, C., and Solorio, T. (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 1298–1308, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1131. URL https://aclanthology.org/N19–1131/.
- Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., and Torralba, A. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis* and Machine Intelligence, 40(6):1452–1464, 2018. doi: 10.1109/TPAMI.2017.2723009.

A. Details on Debiasing Methods

Experiments are conducted on a comprehensive list of commonly used debiasing methods, each of which is designed with special formulation and assumptions.

- Empirical Risk Minimization (ERM) is the standard training method that minimizes the empirical risk on a dataset. This is akin to fine-tuning a pre-trained model on a dataset using cross-entropy loss with no debiasing strategy, which works for both image and text datasets.
- **HypothesisOnly-PoE** (Karimi Mahabadi et al., 2020) assumes the hypothesis part of NLI datasets contains biases. It trains a hypothesis-only (biased) model to measure the bias of each sample, and uses Product-of-Experts (PoE) (Hinton, 2002) to adjust the confidence of the debiased model according to the confidence of the biased model. This approach is evaluated on text datasets.
- WeakLearner-PoE (Sanh et al., 2021) leverages weak learners to capture and model bias, including bias of unknown type. It trains a 2-layer BERT as a biased model and exploits PoE to train the debiased model. This approach is evaluated on text datasets.
- **KernelWhitening** (Gao et al., 2022) aims at learning isotropic sentence embeddings with disentangled robust and spurious representations, with Nyström kernel (Xu et al., 2015). This approach is evaluated on text datasets.
- AttentionPoE (Wang et al., 2023) assumes that the attention to [CLS] token in text classification is biased and introduces PoE on attention weights to learn robust attention patterns for bias mitigation. This approach is evaluated on text datasets.
- σ -Damp (Puli et al., 2023) assuming the standard cross-entropy loss encourages models to prioritize shortcuts over robust features, this model proposes to scale the loss by a temperature. This approach is evaluated on image datasets.
- **DeepFeatReweight** (Kirichenko et al., 2023) discovers that simply retraining the last layer of a neural model–the classification layer in supervised tasks–on top of the existing biased feature extractor is good strategy for bias mitigation. This approach is evaluated on image datasets.
 - **PerSampleGrad** (Ahn et al., 2023) trains a debiased model with non-uniform sampling probability, obtained from per-sample gradient norm of a biased model. This approach is evaluated on image datasets.

B. Implementation details

We follow previous debiasing works for implementation details. For text datasets, we train each debiasing method with Adam optimizer, learning rate 5e - 5, 5 epochs, both KD and Non-KD. For image datasets, we train each debiasing method with Adam optimizer, learning rate 4e - 5, 100 epochs, both KD and Non-KD. For all other hyperparameters, we follow each debiasing method's best-performing setting.

C. Results on Image Datasets

On image datasets, we observe similar results on text datasets. Specifically, we see that KD fall short on distilling the debiasing capabilities. Such ability is transferred more smoothly as teacher and student get similar in scale.

D. Detailed Results on Debiasing Methods and Backbones

We present the detailed results of individual debiasing method and backbone below.

ID OOD Spu. Gap Stu.) $\tau = \tau - 1.2$ 0.5 -1.7 Ò 0.0 2.5 5 -2.5 0.0 T = S - 0.4 - 1.14.5 -3.4 -1.8 0.6 Gap (Tea. *T*=M 0.2 -1.3 -1.3 -0.5 0.3 -4.0 0.0 -1.6 5.4 $\mathcal{T} = B \quad 4.7$ 1.1 -0.1 -0.2 2.5 2.3 -4.3 -2.4 -1.8 -2.5 9.4 5.0 2.2 1.6 -0.2 -0.6 -0.4 4.9 4.5 1.4 0.8 -0.5 -2.9 -2.6 -1.1 -1.4 $\mathcal{T} = L$ 0.1 S=T S=S S=M S=B S=LS = T S = S S = M S = B S = LS=T S=S S=M S=B S=L

Figure 11. **C1: Teacher vs. Student**: average performance gaps between teacher and student models on ID, OOD, and Spurious Gap across image datasets. X-axis and Y-axis show the scale of student (S) and teacher (T) respectively. Each cell shows the performance gap between corresponding scales of a teacher and a student.



Figure 12. C2: Non-KD vs. KD: average performance gaps between Non-KD and KD models on ID, OOD, and Spurious Gap across image datasets. X-axis and Y-axis show the scale of student (S) and teacher (T) respectively. Each cell shows the performance gap between corresponding scales of a teacher and a student.



Figure 13. C2: KD vs. Non-KD: Centralized Kernel Alignment. Highers values indicate higher similarity. X-axis and Y-axis refer to the layers of KD (S) and Non-KD (f_S) respectively.

806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824

803

804

805

770

772 773

774

777

779 780

781

782 783

784

785



Figure 15. C1: Teacher vs. Student: average performance gaps between teacher and student models on ID, OOD, and Spurious Gap on T5.



Figure 17. **C1: Teacher vs. Student**: average performance gaps between teacher and student models on ID, OOD, and Spurious Gap on ViT.





Figure 18. C1: Teacher vs. Student: prediction agreement on BERT. Left: varying S (X-axis) given a fixed teacher with T = L. Right: varying \mathcal{T} (X-axis) given a fixed student with \mathcal{T} = T. Agreement increases as the scale of teacher and student get closer.