

# Think Twice: Measuring the Efficiency of Eliminating Prediction Shortcuts of Question Answering Models

Anonymous ACL submission

## Abstract

While the Large Language Models (LLMs) dominate a majority of language understanding tasks, previous work shows that some of these results are supported by modelling spurious correlations of training datasets. Authors commonly assess model robustness by evaluating their models on out-of-distribution (OOD) datasets of the same task, but these datasets might *share* the bias of the training dataset.

We propose a simple method for measuring a scale of models’ reliance on any identified spurious feature and assess the robustness towards a large set of known and newly found prediction biases for various pre-trained models and debiasing methods in Question Answering (QA). We find that the reported OOD gains of debiasing methods can not be explained by mitigated reliance on biased features, suggesting that biases are shared among different QA datasets. We further evidence this by measuring that performance of OOD models depends on bias features *comparably* to the ID model. Our findings motivate future work to refine the reports of LLMs’ robustness to a level of known spurious features.

## 1 Introduction

Unsupervised pre-training and vast parametrization (Devlin et al., 2018; Radford and Narasimhan, 2018) enable Large Language Models (LLMs) to reach close-to-human accuracy on complex downstream tasks such as Natural Language Inference, Sentiment Analysis, or Question Answering. However, previous work shows that these outstanding results can partially be attributed to models’ reliance on non-representative patterns in training data shared with the test set, such as the high lexical intersection of the entailed hypothesis to premise (Tu et al., 2020) in Natural Language Inference (NLI) or the intersection of the question and answer vocabulary (Shinoda et al., 2021) in extractive Question Answering (QA).

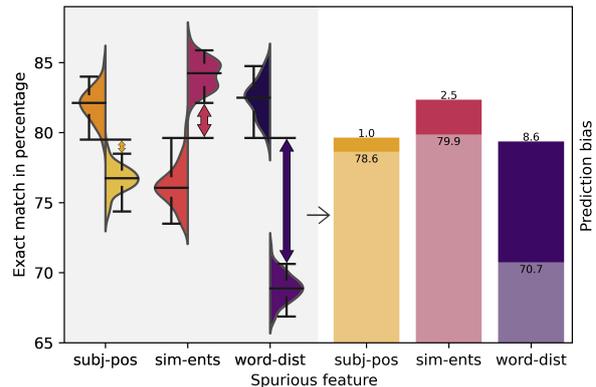


Figure 1: We quantify model reliance on a spurious feature using bootstrapped evaluation on segments of data separated by exploiting chosen bias (left) and subsequently, by measuring the difference in model’s performance over these two groups (right), that we refer to as *Prediction bias* (§3).

A primary motivation for mitigating models’ reliance on such features is to enhance their *robustness* in practice, avoiding fragility to systematic errors when responding the open-ended user requests. Models’ robustness is commonly assessed by measuring prediction quality on samples from other out-of-distribution (OOD) datasets (Clark et al., 2019a; Karimi Mahabadi et al., 2020; Utama et al., 2020b; Xiong et al., 2021). However, the OOD datasets might *share* training biases introduced by shared features, such as data collection methodology, or human annotators’ background (Mehrabi et al., 2021). In such cases, conversely, a model reliant on biased correlations can reach *higher* OOD score despite being more fragile to the adversarial inputs exploiting the biased correlation.

With this motivation, we propose a framework to evaluate models’ reliance on a biased feature in prediction by *splitting* evaluation data to two groups based on a biased feature and *comparing* the prediction quality on these two groups (Fig. 1). This way, we assess a reliance on bias of diverse QA models for several previously and newly identified bias fea-

065 tures identified in this work. Finally, we assess the  
066 efficiency of the state-of-the-art debiasing methods  
067 in mitigating reliance on spurious features over a  
068 resampling baseline and compare the findings to  
069 the commonly-assessed OOD performance.

070 We find that avoiding reliance on spurious fea-  
071 tures does not imply improvements in OOD perfor-  
072 mance; In many cases, debiasing methods mitigate  
073 the model’s prediction bias, but the OOD perfor-  
074 mance drops, while counterintuitively, a magnifi-  
075 cation of bias reliance can also bring large OOD  
076 gains. Aiming to explain this, we directly evaluate  
077 the prediction bias of models trained on different  
078 datasets and confirm that even models trained on  
079 OOD datasets often rely on the *same* spurious cor-  
080 relations as the ID models. This finding motivates the  
081 presented assessment of model robustness towards  
082 known biases, in addition to OOD performance.

083 This paper is structured as follows. Section 2  
084 overviews data biases observed in NLP datasets, re-  
085 cent debiasing methods, and the previous methods  
086 related to measuring inclination to spurious correla-  
087 tions. Section 3 presents our method for measuring  
088 the significance of specific biases. We follow in  
089 Section 4 with details on our evaluation setup, in-  
090 cluding the tested debiasing methods, addressed  
091 bias features, and the design of a set of heuristics  
092 that can exploit them. Subsequently, in Section 5,  
093 we measure and report models’ robustness to bi-  
094 ases and OOD datasets before and after applying  
095 the selected debiasing methods and wrap up our  
096 observations in Sections 6 and 7.

097 **Problem definition** Given a set of inputs  $X =$   
098  $x_{1..i}$  with corresponding labels  $Y = y_{1..i}$  from a  
099 dataset  $\mathcal{D}_{ID}$ , a model  $M$  learns a *task*  $\mathcal{T}$  by identify-  
100 ing *features*  $\mathcal{F}_{1..n}$  that map each  $x_j$  to a correspond-  
101 ing  $y_j$ , assuming that the learned features must be  
102 *consistent* with  $\mathcal{D}_{ID}$ . Since the learned  $\mathcal{F}_{1..n}$  are  
103 distributed in  $M$  and can not be directly evaluated,  
104 we assess whether the learned features are *robust*  
105 for the task  $\mathcal{T}$  by evaluating  $M$  on samples  $X_{OOD}$   
106 of the same task, but drawn from  $\mathcal{D}_{OOD} \not\approx \mathcal{D}_{ID}$ ;  
107 we assume that if  $\mathcal{F}_{1..n} \in M$  are robust, the model  
108 will also perform well on  $X_{OOD}$ . However, the  
109 consistency of the learned  $\mathcal{F}_k$  with both  $X_{ID}$  and  
110  $X_{OOD}$  is merely a necessary and not a sufficient  
111 condition for  $\mathcal{F}_k$  to be robust; If there exists a pair  
112  $(x, y)$  such that the pair is a *valid* sample of the task  
113  $\mathcal{T}$ , but is not consistent with  $\mathcal{F}_k$ , we denote  $\mathcal{F}_k$  as  
114 *spurious* or *bias features* for  $\mathcal{T}$  and refer to models’  
115 reliance on such features as *prediction bias*.

## 2 Background 116

**Spurious correlations of NLP datasets** Previ- 117  
ous work analysing LLMs’ error cases identified 118  
numerous false assumptions that LLMs use in pre- 119  
diction and can be misused to notoriously draw 120  
wrong predictions with the model. 121

122 In Natural Language Inference (NLI), where the 123  
task is to decide whether a pair of sentences entail 124  
one another, McCoy et al. (2019) identifies LLMs’ 125  
reliance on a lexical overlap and on specific shared 126  
syntactic units such as the constituents in the pro- 127  
cessed sentence pair. Asael et al. (2021) identify 128  
the model’s sensitivity to meaning-invariant struc- 129  
ture permutations. Similarly, Chaves and Richter 130  
(2021) identify BERT’s reliance on the invariant 131  
morpho-syntactic composition of the input.

132 In Question Answering, LLMs often rely on the 133  
positional relation of the question and possible an- 134  
swer words, such as falsely assuming their proximi- 135  
ty (Jia and Liang, 2017). Bartolo et al. (2020) find 136  
that models tend to assume that questions and an- 137  
swers contain similar keywords, remaining vulnera- 138  
ble to samples with none or multiple occurrences of 139  
the keywords in the context. Ko et al. (2020) show 140  
models’ preference for the answers in the first two 141  
sentences of the context, being statistically most 142  
likely to answer human-curated questions.

143 A perspective direction circumventing the biases 144  
introduced in data collection is presented in adver- 145  
sarial data collection (Jia and Liang, 2017; Bartolo 146  
et al., 2020) where the annotators collect the dataset 147  
with the intention of fooling the possibly-biased 148  
model, possibly enhancing the model-in-the-loop 149  
in several iterations. Still, some doubts remain; 150  
Models trained on adversarial data may work better 151  
on adversarial datasets but underperform on other 152  
OOD datasets (Kaushik et al., 2021), or introduce 153  
its own set of biases (Kovatchev et al., 2022).

**Debiasing methods** A well-established line of 154  
work proposes to address the known dataset bi- 155  
ases in the training process. Karimi Mahabadi 156  
et al. (2020) and He et al. (2019) obtain the de- 157  
biased model by (i) training a *biased model* that 158  
exploits the unwanted bias, and (ii) training the 159  
debiased model as a complement to the biased one 160  
in a Product-of-Experts (PoE) framework (Hinton, 161  
2002). Clark et al. (2019a) extend this framework 162  
in the LearnedMixin method, learning to weigh 163  
the contribution of the biased and debiased model 164  
in the complementary ensemble. Niu and Zhang 165  
(2021) simulate the model for non-biased, out-of- 166

distribution dataset through counterfactual reasoning (Niu et al., 2021) and use the resulting distribution for distilling target (Hinton et al., 2015), similarly to the LearnedMixin. Biased samples can also be identified in other ways, for instance, by the model’s overconfidence (Wu et al., 2020).

In a complement to PoE approaches, other works apply model confidence regularization on the samples denoted as biased. Feng et al. (2018) and Utama et al. (2020a) downweigh the predicted probability of the examples marked as biased by humans or a model. Xiong et al. (2021) find that a more precise calibration of the biased model might bring further benefits to this framework, consistently to our observations. Distributionally Robust Optimization (DRO) methods are another group of reweighting algorithms, addressing assumed imperfection of training datasets by (i) segmenting data into groups of diverse covariate shifts (Sagawa et al., 2020) and (ii) minimizing the worst-case risk over all groups (Zhou et al., 2021). We note that our bias measurement method closely relates to group DRO methods and can, for instance, also serve as a method for quantifying per-group risk.

**Robustness measures** Most of the work on enhancing models’ robustness evaluates the acquired robustness on OOD datasets. In some cases, the evaluation utilizes datasets specially constructed to exploit the biases typical for a given task, such as HANS (McCoy et al., 2019) for NLI, PAWS (Zhang et al., 2019) for Paraphrase Identification, or AdversarialQA (Bartolo et al., 2020) for Question Answering, that we also use in evaluations.

Similar to us, some previous work quantified dataset biases by splitting data into two subsets, comparing model behaviour between these groups. McCoy et al. (2019) perform such evaluation over MNLI, demonstrating large margins in accuracy over the two groups and superior robustness of BERT over previous models. Similarly, Utama et al. (2020b) compare two groups based on prediction confidence. Our Prediction bias measure follows a similar approach in QA but provides a more reliable assessment thanks to bootstrapping. Compared to the previous work, we assess models’ reliance on a range of 7 spurious features, making overall conclusions more robust.

An ability to measure a model’s reliance on undesired features is well-applicable in quantifying socially problematic biases. Previous work also utilizes specialized domain knowledge in models’

```

func measure_bias(M, X, h, Th):
    Ah ← h(X)
    X1 ← x1 ∈ X : Ah(x1) ≤ Th
    X2 ← x2 ∈ X : Ah(x2) > Th
    foreach X1l ∈ repeat(sample(X1)) do
        | E1 ← E1 + evaluate(M(X1l))
    foreach X2l ∈ repeat(sample(X2)) do
        | E2 ← E2 + evaluate(M(X2l))
    dist ← max(0; E1↓ − E2↑; E2↓ − E1↑)
    return dist

```

**Algorithm 1:** We measure *Prediction bias* of the model *M* exploited by the *heuristic h* on dataset *X*, as a *difference* of *M*’s performance on two groups (*X*<sub>1</sub> and *X*<sub>2</sub>) obtained by segmenting the samples of *X* by the *attribute A<sub>h</sub> = h(X)* on a given threshold *T<sub>h</sub>*.

We bootstrap both evaluations, (*samples* = 800, *trials* = 100, and obtain two sets of measurements (*E*<sub>1</sub> and *E*<sub>2</sub>), of which we subtract the upper and lower quantiles *E*<sup>↑</sup> and *E*<sup>↓</sup> (*q*<sup>↑</sup> = 0.975, *q*<sup>↓</sup> = 0.025) and consider such distance a scale of the learned prediction bias.

bias evaluation but might not scale to other bias features; Parrish et al. (2022) collect ambiguous contexts and assess the models’ inclination to utilize stereotypes as prediction features. Bordia and Bowman (2019) quantify LM’s gender bias by the co-occurrence of selected gender-associated words with gender-ambiguous words, such as *doctor*.

### 3 Measuring Prediction Bias

We assess a model’s sensitivity to a known spurious feature in the following sequence of steps. This methodology is also visualized in Figure 1 and described in Algorithm 1.

We start by (i) implementing a *heuristic*, i.e. a method *h* : *X* → ℝ, that for *all* samples of dataset *X* computes an *attribute A<sub>h</sub> ∈ ℝ* corresponding to the feature *F* that we susprise as non-representative, yet predictive for our training set and (ii) evaluate *h* on evaluation dataset *X*. (iii) We choose a threshold *T<sub>h</sub>* that we use to (iv) split the dataset into two segments by *A<sub>h</sub>*. Finally, (v) we evaluate the assessed model *M* on both of these segments, in our case using Exact match measure, and (vi) measure model **prediction bias** as the *difference* in performance between these two groups. Using bootstrapped evaluation, we mitigate the effect of randomness by only comparing selected

quantiles of confidence intervals. We propose to perform a hyperparameter search for the heuristic’s threshold  $T_h$  that maximizes the measured distance.

**Interpretation** Given the reliance on bootstrapping, we state that model’s *true* performance polarisation is  $0.975 \times 0.975 = 95.06\%$ -likely to be equal or higher than the measured Prediction bias (with  $q^\uparrow = 0.975, q^\downarrow = 0.025$  as in Algorithm 1).

Nevertheless, one should note that the proposed measure should not be used in a standalone but rather in a complement to an ID evaluation, as one can reduce the Prediction bias merely by *lowering* the performance on the better-performing ID subset. Therefore, we report the values of Prediction bias together with the performance on a worse-performing, i.e. presumably non-biased split.

Another consideration concerns the “natural” polarisation of difficulty between samples; That is a portion of Prediction bias which can be explained by the features  $\mathcal{F}$  that are *representative* for the evaluated task (§1). One should note that the reduction of Prediction bias is meaningful only up to the level of the natural sample difficulty.

The validation set of SQuAD contains the annotations by three annotators that we use to quantify a level of Prediction bias that can be explained by the questions’ natural difficulty (further denoted as *Human* model); We report the minimum over Prediction biases of the annotators among each other.

Finally, even though we perform a hyperparameter search for optimal heuristics’ thresholds  $T_h$  feasible for a given size of dataset splits, there are no guarantees on the maximality of the found  $T_h$ . Hence, Prediction bias only provides the *lower bounds* of the model’s worst-case polarisation.

## 4 Experiments

Our main objective is to assess the efficiency of different training decisions in mitigating the reliance of the model on spurious correlations that we assume to be present in a dataset. In QA, we identify the spurious covariates in SQuAD dataset (Rajpurkar et al., 2016), with existing work documenting a variety of learnt spurious correlations.

For each suspected bias feature, we first describe and implement the exploiting heuristics that we use to segment groups in the Prediction bias measure (§4.1). Subsequently, we observe the impact of the selected pre-training strategies (§4.2) and debiasing methods designed to address the over-reliance on

biased features (§4.3 – §4.4) on the Prediction bias and OOD performance of the resulting models.

### 4.1 Biases and Exploiting Heuristics

Our work extends the list of previously-reported QA biases based on our experience with two novel bias features that we later assess as significant. The spurious features newly identified in this work are preceded with  $\pm$ .

Together with each bias, we also briefly describe its exploiting heuristic computing the non-representative feature  $A_h$  (Algorithm 1).

**Distance of Question words from Answer words (*word-dist*)** Jia and Liang (2017) propose that the models are prone to return answers close to the vocabulary of the question in context. Hence, *word-dist* computes how close the closest question word is to the first answer in the context and computes the distance ( $A_h$ ) as a number of words between the closest question word and the answer span.

**Similar words between Question and Context (*sim-word*)** Shinoda et al. (2021) report the common occurrence of a high lexical overlap between the question and the correct answer over QA datasets. In *sim-word* heuristic, we represent the lexical overlap by the number of shared words between the question and the context. Both are defined as sets, and the intersection size of these two sets is computed as the heuristic’s evaluation ( $A_h$ ).

**Answer position in Context (*ans-pos*)** Ko et al. (2020) report that QA models may learn to falsely assume the answer’s occurrence in the first two sentences. The exploiting heuristic first segments the context into sentences, then identifies the sentence containing the answer and yields a scalar corresponding to the rank of the sentence within the context that contains the answer ( $A_h$ ).

**Cosine similarity of Question and Answer (*cos-sim*)** Clark et al. (2019a) use the TF-IDF similarity as a biased model for QA, implicitly identifying a bias in undesired reliance of the model on the match of the keywords between the question and retrieved answer. We exploit this feature by (i) fitting the TF-IDF model on all SQuAD contexts, (ii) inferring the TF-IDF vectors of both questions and their corresponding answers, and (iii) returning the scalar ( $A_h$ ) as cosine similarity between the TF-IDF vectors of question and answer.

**Answer length (*ans-len*)** Bartolo et al. (2020) show that QA models trained on SQuAD make errors much more often on questions asking for

longer answers, implicitly identifying models' reliance on a feature that the answer must comprise at most a few words. We exploit this feature by simply computing  $A_h$  as the length of the answer.

**+Number of Question's Named Entities in Context (*sim-ents*)** We suspect that the in-context presence of multiple named entities, such as multiple personal names or locations, might perplex the QA model's prediction. This might suggest that models tend to reduce the QA task to a simpler yet irrelevant problem of Named Entity Recognition. We utilize a pre-trained BERT NER model provided within SPACY library (Honnibal and Montani, 2017) to identify named entities of the *question type* (i.e., *personal names* if the question starts with "Who"). Then, we count  $A_h$  as the number of matching named entities in the context.

**+Position of Question's subject to the correct Answer in Context (*subj-pos*)** Our observations suggest that the position of the question's subject in the context impacts the predicted answer spans of QA models. In the corresponding heuristic, using SPACY library, we (i) identify the questions' subject expression and (ii) locate its occurrences in the context. We (iii) locate the answer span and compute  $A_h$  as a relative position of the answer: either before the subject, after the subject, or after multiple occurrences of the question subject.

## 4.2 Evaluated Models

To estimate the impact of selected pre-training strategies on the robustness of the resulting model, we conventionally fine-tune a set of diverse pre-trained LLMs for extractive QA.

We alternate between the following models: BERT-BASE (Devlin et al., 2019), ROBERTA-BASE and ROBERTA-LARGE (Liu et al., 2019), ELECTRA-BASE (Clark et al., 2020) and T5-LARGE (Raffel et al., 2020). This selection allows us to outline the impact of the various features on the robustness of the final QA model: (i) pre-training data volume (BERT-BASE vs ROBERTA-BASE), (ii) model size (ROBERTA-BASE vs ROBERTA-LARGE), (iii) pre-training objective (BERT-BASE vs ELECTRA-BASE), or (iv) extractive vs. generative prediction mode (T5 vs. others).

We also evaluate the prediction bias of recent multi-task in-context learners, without fine-tuning: T0 (Sanh et al., 2022) trained for zero-shot in-context learning excluding SQuAD, and FLAN-T5 (Chung et al., 2022) trained on a mixture of more than 1,800 tasks, including SQuAD.

## 4.3 Debiasing Baseline: Resampling (RESAM)

Based on the heuristics and their tuned configuration, our baseline method performs simple super-sampling of the underrepresented group ( $X_1$  or  $X_2$  in Algorithm 1) until the two groups are represented equally. This approach shows the possibility of bias reduction by simply normalizing the distribution of the biased samples in the dataset, requiring only the identification of the members of the under-represented group. RESAM closely follows the routine of Algorithm 1 and splits the data by the optimal threshold of the attributes of the heuristics corresponding to each addressed bias.

## 4.4 Assessed Debiasing Methods

We assess the efficiency of debiasing methods in eliminating Prediction bias for the representatives of two diverse debiasing methods. In addition to Prediction bias, we also report the resulting performance on three OOD datasets. We follow the reference implementations as closely as possible while scaling the scope of experiments from one to seven separately-addressed biases. Complete description of training settings is in Appendix B.2.

**LearnedMixin (LMIX)** method (Clark et al., 2019b) is a popular adaptation of Product-of-Experts framework (Hinton, 2002), with a set of refinements (§2), that uses a *biased model* as a complement of the trained debiased model in a weighted composition. We reimplement the reference implementation with the following alterations. Instead of the BIDAf model, we use stronger BERT-BASE as the trained debiased model. Instead of using a TF-IDF-based bias model customized for a single bias type, we opt for a universal approach for obtaining biased models (Appendix B.2.1). We rerun the parameter search and use a different entropy penalty ( $H = 0.4$ ) throughout all experiments.

**Confidence Regularization (CREG)** aims to reduce the model's confidence, i.e. the predicted score over samples marked as biased. Utama et al. (2020a) propose to reduce the confidence of the biased samples using a distillation from the conventional QA teacher model, scaled down by the relative scores of a biased predictor. In our experiments, we consistently use BERT-BASE for both the teacher and bias model. To enable comparability with LMIX, we use identical bias models for both methods (Described in Appendix B.2.1).

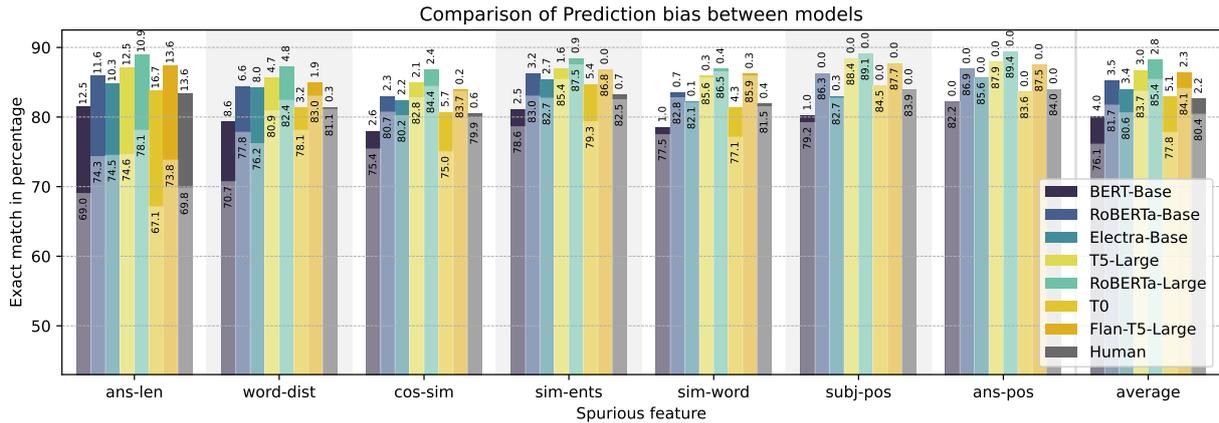


Figure 2: **Prediction bias per pre-trained model.** The worse-performing split performance (lower bars) and Prediction bias (upper bars, sorted by group average) of QA models trained from different pre-trained LLMs, trained and evaluated on SQuAD for Exact match. Per-group bootstrapping of 100 repeats with 800 samples.

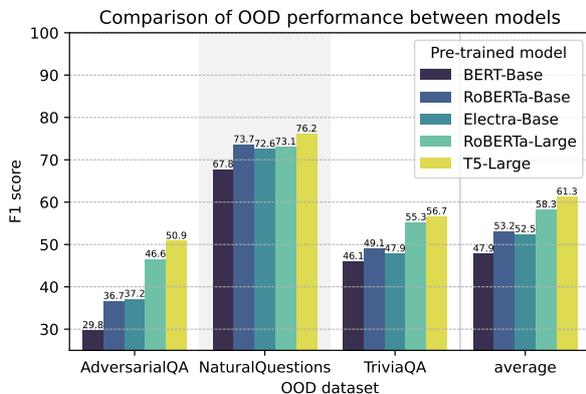


Figure 3: **OOD performance per pre-trained model.** Comparison of F1-score of different models fine-tuned on SQuAD and evaluated on listed OOD datasets.

## 5 Results

### 5.1 Impact of Pre-training

Figure 2 compares the Prediction bias of the fine-tuned models of diverse pre-training data volumes and objectives, followed by in-context learning models and human reference.

The results suggest that increased amounts of pre-training data of the base models (cf. BERT-BASE and others) might mitigate the models’ reliance on the bias. The results are less conclusive in a comparison of different pre-training objectives (cf. ROBERTA-BASE and ELECTRA-BASE); While ELECTRA is less polarised in 4 out of 7 cases, the differences are minimal. The largest reduction of Prediction bias ( $-1.2$  on average) is achieved by increasing the model size of ROBERTA-LARGE.

Analogously, Figure 3 compares OOD performance on selected QA datasets: AdversarialQA (Jia and Liang, 2017), NaturalQuestions (Kwiatkowski et al., 2019) and TriviaQA (Joshi

et al., 2017). The concluding robustness ranking is mainly consistent with the Prediction bias ranking, with an exception of generative fine-tuning (T5), which outperforms others on OOD datasets but not on a reduction of the reliance on spurious features.

### 5.2 Prediction bias of OOD models

Figure 4 compares Prediction bias over the least-biased ROBERTA-LARGE models trained on different datasets. All evaluations are split on heuristics’ thresholds  $T_h$  optimal for SQuAD model, which allows comparability to the shared human reference but implies that larger Prediction bias for OOD models might exist. We see that all Prediction biases learnt on SQuAD are also learnt from at least one OOD dataset. For Trivia model, *all* types of biases identified in SQuAD are magnified.

We specifically note the comparison of Prediction bias of the SQuAD model to the model trained on AdversarialQA, collected adversarially to a SQuAD model; We find that AdversarialQA model is the only OOD model lowering reliance on all biased features that are over the level of natural bias, supporting the argued efficiency of adversarial data collection in addressing original dataset biases.

### 5.3 Impact of Debiasing

Figure 5 compares the biases of Question Answering models obtained within three debiasing methods (§4.3 – §4.4), applied to the most-biased BERT-BASE model. We observe that debiasing methods are not consistent in the efficiency of mitigating the reliance on the addressed bias feature. In fact, only RESAM baseline lowers the bias of the original model consistently. We attribute this inconsistency to methods’ sensitivity to *bias model*,

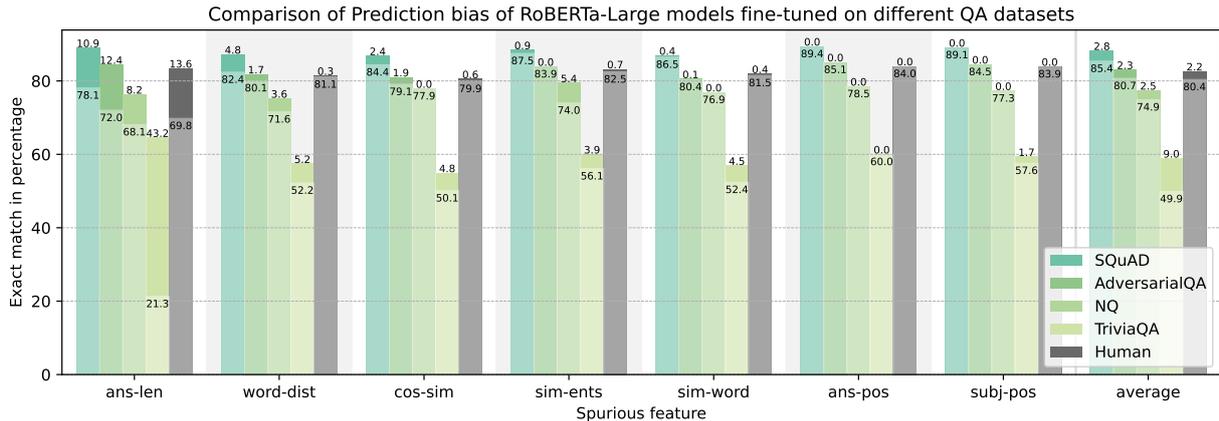


Figure 4: **Prediction bias per dataset.** The worse-performing split performance (lower bars) and Prediction bias (upper bars) of ROBERTA-LARGE trained on different QA datasets, evaluated on a validation split of SQuAD for Exact match. All evaluation splits are identical, identified as maximal for the SQuAD-trained model (Appx. C).

Table 1: **OOD performance of debiasing methods.** Differences of F1-scores of QA models trained on SQuAD using specified debiasing methods (§4.4) to address selected bias features (§4.1) evaluated on three OOD datasets; *AdversarialQA* / *NaturalQuestions* / *TriviaQA*, respectively. Largest gains per dataset are **bold**.

	Original model			29.8 / 67.8 / 46.1		
	ReSam		LMix	CReg		
	<i>AQA</i> / <i>NQ</i> / <i>Trivia</i>					
<i>ans-len</i>	-0.8 / -5.6 / -1.7	-0.9 / -19.7 / -3.3	-0.4 / +5.5 / +2.1			
<i>word-dist</i>	+0.5 / +1.3 / +0.0	+0.9 / -6.4 / +1.5	<b>+1.4</b> / <b>+7.5</b> / -0.5			
<i>cos-sim</i>	-0.1 / +0.3 / -1.3	+0.4 / -11.3 / -4.1	-0.3 / +7.4 / +1.1			
<i>sim-ents</i>	+1.1 / +1.5 / +0.3	-0.1 / -9.5 / -1.2	-1.0 / +5.9 / +2.0			
<i>sim-word</i>	+0.3 / +0.1 / +0.4	-0.3 / -21.4 / -2.9	-0.7 / +3.9 / +1.4			
<i>subj-pos</i>	-1.6 / -0.7 / -2.2	-1.3 / -14.8 / -1.3	+0.0 / +5.1 / +1.6			
<i>Average</i>	-0.45	-5.31	+2.33			

further discussed in §6. While LMIX is the most efficient in addressing Prediction bias in average, consistently to Clark et al. (2019a) we see that this often comes for a price of the ID performance.

Table 1 enumerates the OOD performance of debiased models over three diverse QA datasets. By comparing these results to Prediction bias (Fig. 5), we see many cases where the reduction of Prediction bias can not explain improvements of OOD; For instance, addressing *word-dist* bias using CREG improves OOD performance by 2.8% of an exact match on average and by 7.5 on *NaturalQuestions*, but the Prediction bias of such model increases by 1.1 points. Similarly, CREG addressing *sim-word* bias, delivers 1.5-point average gain on OOD but raises Prediction bias by 0.9 points.

Figure 6 further evaluates the impact of addressing one bias to other known biases in cases where each method delivers the largest Prediction bias reduction. We see that addressing a specific bias also affects the scope of the model’s reliance on other covariates. Results suggest that CREG might be

more robust to enlarging of other biases, increasing other Prediction biases by 0.31 on average, as compared to LMIX (0.6) and RESAM (0.38).

## 6 Discussion

**Pre-training and models’ robustness** The bias-level analyses of diverse pre-trained models (§5.1) suggest that the mere increase of pre-training data and model parameters guide the fine-tuned models to lower reliance on biased features. However, we can find exceptions, such as in the case of ROBERTA-LARGE and ELECTRA-BASE on *ans-len*. We speculate that even larger volumes of data might make the model more attracted to taking a shortcut through easier problem formulations, such as through Named entity recognition (cf. BERT-BASE and ROBERTA-BASE on *sim-ents*).

Comparing the prediction bias of in-context learners with the fine-tuned models, we see that multi-task learning does not necessarily result in lower prediction bias or increased performance in the harder group; While FLAN-T5 on average reduces bias almost to the human level, T0’s quality is affected by spurious features even more than the models fine-tuned on biased SQuAD.

### OOD performance and Prediction bias relation

Our results conclude that the previously-reported improvements in OOD performance attributed to the debiasing might not be attributed to the mitigated reliance on a spurious correlation; (i) We measure that Prediction bias of the models trained directly on OOD datasets is still present over the level of human Prediction bias (§5.2). Therefore, it is possible to maintain OOD gains by learning to rely on bias features. (ii) In practice, we find cases

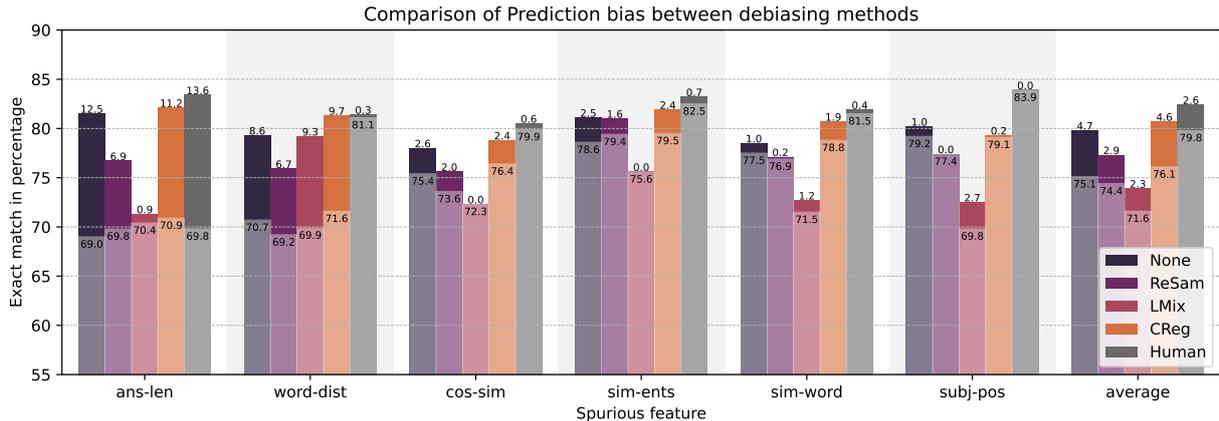


Figure 5: **Prediction bias per debiasing methods.** The worse-performing split performance (lower bars) and Prediction bias (upper bars) of BERT-BASE trained using selected debiasing methods, evaluated for Exact match on validation SQuAD. Per-group evaluations were measured using bootstrapping of 100 repeats with 800 samples.

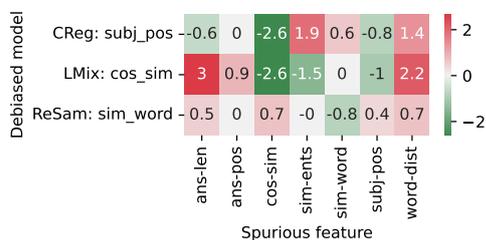


Figure 6: **Cross-bias evaluation of debiased models.** A relative change of Prediction bias by all spurious correlations, caused by applying inspected debiasing methods on BERT-BASE QA model, in addressing specified spurious correlation. A full matrix is in Appx. A, Fig. 7.

where applying a debiasing method *magnifies* Prediction bias, but the resulting model still performs better in most OOD evaluations (§5.3).

### Practical aspects of applying debiasing methods

While we confirm that debiasing methods enable improvements in the OOD, we find that the significance of such improvements largely varies between the addressed biases and the suitable configuration for one bias and dataset pair is often suboptimal for others. The scope of this variance can be seen in Table 1 from the comparison of average OOD performance of LMIX and CREG on *word-dist*, used to pick methods’ hyperparameters and bias models (Appendix B.2), and other biases; Both of the methods perform best on the bias used in parameter tuning, and the differences are often large. Bias-specific parameter tuning is further convoluted by the speed of the convergence of debiasing methods, which we measure as approximately 4-times slower for CREG and 3.5-times slower for LMIX, compared to the standard fine-tuning of QA models.

The bias model is an important parameter of both assessed debiasing methods. We find that the

scores have to be rescaled for trained bias models to avoid perplexing the trained model on biased samples and that the optimal scaling parameter is also bias-specific. The selection of the bias model also affects the optimal Entropy scaling  $H$  of LMIX; we find that the optimal value ( $H = 2.0$ ) for AdversarialQA reported by LMIX authors is also not close to optimal ( $H = 0.4$ ) for our bias model.

## 7 Conclusion

Our work sets out to investigate the impact of various training decisions, including different pre-training and debiasing strategies, on models’ reliance on specific spurious features in QA, complementing the commonly-used out-of-distribution evaluations. We use SQuAD to survey documented and identify some new biased features but evaluate the reliance on these features for models trained on four different QA datasets.

We find that (i) the OOD performance of different base models usually corresponds to models’ reliance on bias features. However, (ii) the state-of-the-art debiasing methods can improve OOD performance *without* minimising the model’s reliance on spurious features, suggesting that dataset biases might be *shared* among QA datasets. (iii) We further evidence this by measuring the reliance on a spurious feature of models trained on other (OOD) datasets and find OOD models similarly or even more *reliant* on spurious features of SQuAD.

These findings aim to motivate future work to assess models’ robustness also on a level of specific bias features, evading false conclusions on models’ robustness, ultimately fostering progress toward reliable and socially unbiased language models.

## References

- 610  
611 Dimion Asael, Zachary Ziegler, and Yonatan Belinkov. 2021. [A generative approach for mitigating structural](#)  
612 [biases in natural language inference](#). *arXiv preprint*  
613 *arXiv:2108.14006*. 667  
614
- 615 Max Bartolo, A Roberts, Johannes Welbl, Sebastian  
616 Riedel, and Pontus Stenetorp. 2020. [Beat the ai: In-](#)  
617 [vestigating adversarial human annotation for reading](#)  
618 [comprehension](#). *Transactions of the Association for*  
619 *Computational Linguistics*, 8:662–678. 668
- 620 Shikha Bordia and Samuel R. Bowman. 2019. [Identify-](#)  
621 [ing and reducing gender bias in word-level language](#)  
622 [models](#). In *Proceedings of the 2019 Conference of the*  
623 *North American Chapter of the Association for Com-*  
624 *putational Linguistics: Student Research Workshop*,  
625 pages 7–15, Minneapolis, Minnesota. Association for  
626 Computational Linguistics. 669
- 627 Rui P. Chaves and Stephanie N. Richter. 2021. [Look](#)  
628 [at that! BERT can be easily distracted from paying](#)  
629 [attention to morphosyntax](#). In *Proceedings of the*  
630 *Society for Computation in Linguistics 2021*, pages  
631 28–38, Online. Association for Computational Lin-  
632 guistics. 670
- 633 Hyung Won Chung, Le Hou, Shayne Longpre, Barret  
634 Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi  
635 Wang, Mostafa Dehghani, Siddhartha Brahma, Al-  
636 bert Webson, Shixiang Shane Gu, Zhuyun Dai,  
637 Mirac Suzgun, Xinyun Chen, Aakanksha Chowd-  
638 hery, Alex Castro-Ros, Marie Pellat, Kevin Robin-  
639 son, Dasha Valter, Sharan Narang, Gaurav Mishra,  
640 Adams Yu, Vincent Zhao, Yanping Huang, Andrew  
641 Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean,  
642 Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V.  
643 Le, and Jason Wei. 2022. [Scaling Instruction-](#)  
644 [Finetuned Language Models](#). *arXiv e-prints*, page  
645 arXiv:2210.11416. 671
- 646 Christopher Clark, Mark Yatskar, and Luke Zettlemoyer.  
647 2019a. [Don’t take the easy way out: Ensemble](#)  
648 [based methods for avoiding known dataset biases](#).  
649 In *Proceedings of the 2019 Conference on Empirical*  
650 *Methods in Natural Language Processing and the*  
651 *9th International Joint Conference on Natural Lan-*  
652 *guage Processing (EMNLP-IJCNLP)*, pages 4069–  
653 4082, Hong Kong, China. Association for Computa-  
654 tional Linguistics. 672
- 655 Kevin Clark, Urvashi Khandelwal, Omer Levy, and  
656 Christopher D. Manning. 2019b. [What does BERT](#)  
657 [look at? an analysis of BERT’s attention](#). In *Proc.*  
658 *of the 2019 ACL Workshop BlackboxNLP: Analyzing*  
659 *and Interpreting Neural Networks for NLP*, pages  
660 276–286, Florence, Italy. ACL. 673
- 661 Kevin Clark, Minh-Thang Luong, Quoc V. Le, and  
662 Christopher D. Manning. 2020. [ELECTRA: Pre-](#)  
663 [training Text Encoders as Discriminators Rather](#)  
664 [Than Generators](#). *CoRR*, abs/2003.10555v1. 674
- 665 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
666 Kristina Toutanova. 2018. [BERT: Pre-training of](#)  
[deep bidirectional transformers for language under-](#)  
[standing](#). *CoRR*, abs/1810.04805v2. 675
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
Kristina Toutanova. 2019. [BERT: Pre-training of](#)  
[Deep Bidirectional Transformers for Language Un-](#)  
[derstanding](#). In *Proc. of the 2019 Conference of*  
*the NAACL: Human Language Technologies*, pages  
4171–4186, Minneapolis, USA. ACL. 676
- Shi Feng, Eric Wallace, Alvin Grissom II, Mohit Iyyer,  
Pedro Rodriguez, and Jordan Boyd-Graber. 2018.  
[Pathologies of neural models make interpretations](#)  
[difficult](#). In *Proceedings of the 2018 Conference on*  
*Empirical Methods in Natural Language Processing*,  
pages 3719–3728, Brussels, Belgium. Association  
for Computational Linguistics. 677
- He He, Sheng Zha, and Haohan Wang. 2019. [Unlearn](#)  
[dataset bias in natural language inference by fitting](#)  
[the residual](#). In *Proceedings of the 2nd Workshop on*  
*Deep Learning Approaches for Low-Resource NLP*  
*(DeepLo 2019)*, pages 132–142, Hong Kong, China.  
Association for Computational Linguistics. 678
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015.  
[Distilling the knowledge in a neural network](#). Cite  
arxiv:1503.02531Comment: NIPS 2014 Deep Learn-  
ing Workshop. 679
- Geoffrey E. Hinton. 2002. [Training Products of Ex-](#)  
[perts by Minimizing Contrastive Divergence](#). *Neural*  
*Computation*, 14(8):1771–1800. 680
- Matthew Honnibal and Ines Montani. 2017. [spaCy 2:](#)  
[Natural language understanding with Bloom embed-](#)  
[dings, convolutional neural networks and incremental](#)  
[parsing](#). 681
- Robin Jia and Percy Liang. 2017. [Adversarial exam-](#)  
[ples for evaluating reading comprehension systems](#).  
In *Proceedings of the 2017 Conference on Empiri-*  
*cal Methods in Natural Language Processing*, pages  
2021–2031, Copenhagen, Denmark. Association for  
Computational Linguistics. 682
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke  
Zettlemoyer. 2017. [Triviaqa: A large scale distantly](#)  
[supervised challenge dataset for reading comprehen-](#)  
[sion](#). *arXiv preprint arXiv:1705.03551*. 683
- Rabeeh Karimi Mahabadi, Yonatan Belinkov, and James  
Henderson. 2020. [End-to-end bias mitigation by](#)  
[modelling biases in corpora](#). In *Proceedings of the*  
*58th Annual Meeting of the Association for Compu-*  
*tational Linguistics*, pages 8706–8716, Online. Asso-  
ciation for Computational Linguistics. 684
- Divyansh Kaushik, Douwe Kiela, Zachary C. Lipton,  
and Wen-tau Yih. 2021. [On the efficacy of adversar-](#)  
[ial data collection for question answering: Results](#)  
[from a large-scale randomized study](#). In *Proceedings*  
*of the 59th Annual Meeting of the Association for*  
*Computational Linguistics and the 11th International*  
*Joint Conference on Natural Language Processing*  
*(Volume 1: Long Papers)*, pages 6618–6633, Online.  
Association for Computational Linguistics. 685



839 Thomas Wolf, Julien Chaumond, Lysandre Debut, Vic-  
 840 tor Sanh, Clement Delangue, Anthony Moi, Pier-  
 841 ric Cistac, Morgan Funtowicz, Joe Davison, Sam  
 842 Shleifer, et al. 2020a. **Transformers: State-of-the-**  
 843 **art natural language processing.** In *Proceedings of*  
 844 *the 2020 Conference on Empirical Methods in Natu-*  
 845 *ral Language Processing: System Demonstrations,*  
 846 pages 38–45.

847 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien  
 848 Chaumond, Clement Delangue, Anthony Moi, Pier-  
 849 ric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz,  
 850 Joe Davison, Sam Shleifer, Patrick von Platen, Clara  
 851 Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven  
 852 Le Scao, Sylvain Gugger, Mariama Drame, Quentin  
 853 Lhoest, and Alexander Rush. 2020b. **Transformers:**  
 854 **State-of-the-Art Natural Language Processing.** In  
 855 *Proc. of the 2020 Conf. EMNLP: System Demonstrations,*  
 856 pages 38–45. ACL.

857 Mingzhu Wu, Nafise Sadat Moosavi, Andreas Rücklé,  
 858 and Iryna Gurevych. 2020. **Improving QA general-**  
 859 **ization by concurrent modeling of multiple biases.** In  
 860 *Findings of the Association for Computational Lin-*  
 861 *guistics: EMNLP 2020,* pages 839–853. Association  
 862 for Computational Linguistics.

863 Ruibin Xiong, Yimeng Chen, Liang Pang, Xueqi Cheng,  
 864 Zhi-Ming Ma, and Yanyan Lan. 2021. **Uncertainty**  
 865 **calibration for ensemble-based debiasing methods.**  
 866 In *Advances in Neural Information Processing Sys-*  
 867 *tems.*

868 Yuan Zhang, Jason Baldridge, and Luheng He. 2019.  
 869 **PAWS: Paraphrase Adversaries from Word Scram-**  
 870 **bling.** In *Proc. of the 2019 Conf. NAACL-HLT,* pages  
 871 1298–1308, Minneapolis, USA. ACL.

872 Chunting Zhou, Xuezhe Ma, Paul Michel, and Graham  
 873 Neubig. 2021. **Examining and combating spurious**  
 874 **features under distribution shift.** In *Proceedings of*  
 875 *the 38th International Conference on Machine Learn-*  
 876 *ing,* volume 139 of *Proceedings of Machine Learning*  
 877 *Research,* pages 12857–12867. PMLR.

## 878 A Cross-Bias Matrix of All Debaised 879 Models

880 Figure 7 shows the change of Prediction bias by  
 881 applying the listed debiasing methods to eliminate  
 882 the associated bias feature. We see that some biases  
 883 are more difficult to address, while other ones can  
 884 be transitively addressed through others.

## 885 B Details of Training Configurations

886 This section overviews all configurations that we  
 887 have set in training the debaised models (§4.3 – 4.4)  
 888 as well as the conventional QA fine-tuning compar-  
 889 ing the impact of pre-training on QA models’  
 890 robustness (§4.2).

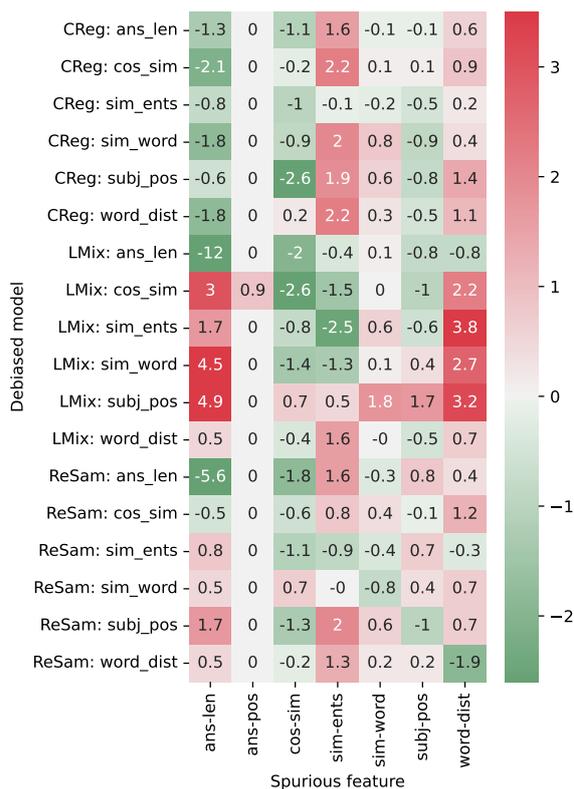


Figure 7: **Full cross-bias evaluation of debaised mod-**  
**els.** A relative change of Prediction bias by all spurious  
 correlations, caused by applying inspected debiasing  
 methods on BERT-BASE QA model, in addressing spec-  
 ified spurious correlation.

## 891 B.1 Standard Fine-tuning

892 For model fine-tuning, we use following hyperpa-  
 893 rameters: **learning rate:**  $2e^{-5}$ , **batch size:** 16,  
 894 **evaluation:** each 200 steps and **train epochs:** 3.  
 895 We also set the **early stopping patience** to 10 eval-  
 896 uation steps, based on a validation loss of the train-  
 897 ing dataset (SQuAD) also used for selecting the  
 898 evaluated model. The **validation loss** of the eval-  
 899 uated model is 1.02. All other parameters can be  
 900 retrieved from the defaults of TrainingArguments  
 901 of HuggingFace (Wolf et al., 2020b) in version  
 902 4.19.1.

903 We use the listed configuration also in training  
 904 the generative T5 model. We use the Adaptor li-  
 905 brary (Štefánik et al., 2022) in version 0.1.6 for  
 906 fine-tuning T5 for generating answers.

## 907 B.2 Debiasing Training Experiments

### 908 B.2.1 Bias models

909 The canonical debiasing implementations utilize  
 910 bias-specific models for identifying bias; Clark  
 911 et al. (2019b) use the TF-IDF model as a scalar

of possible bias for each QA sample, while Utama et al. (2020a) experiment with a percentage of the shared words and cosine embeddings between word distances, in NLI context.

As we scale our experiments to six different biases, we opt for a universal approach for obtaining bias models for both LMIX and CREG and train each bias’ model on a better-performing segment of the dataset identified using the approach described in Section 3. For all our biased models, we train BERT-BASE architecture from scratch and pick the checkpoint with a maximal difference of the F1-score between the two segments from the validation split of SQuAD.

While our approach scales well over many biases, a significant difference between the learned bias models original ones, such as TF-IDF, is the *scale* of prediction probabilities; As the trained bias models become very confident on a biased subset, often reaching probabilities close to 1 for the biased samples. A “perfect” bias model causes problems for both LMIX and CREG as such model forces the trained model to avoid correct predictions on the biased samples completely. We learn to address this problem by rescaling bias predictions and tuning the scaling interval based on a validation performance of the debiased model. Consequently, we scale the bias probabilities to  $\langle 0; 0.2 \rangle$  for LMIX and  $\langle 0; 0.1 \rangle$  for CREG. Further details on bias models can be found in Appendix B.2.

In the initial phase, we experiment with diverse configurations and sizes of bias models, intending to maximize the polarization of performance on the biased and non-biased subsets. Among different configurations of model sizes and configurations, we find that the highest polarisation can be reached using BERT-BASE architecture trained from scratch. We fix this decision and the parameters (learning rate  $4e^{-5}$ , a number of training steps 88,000) with respect to the maximum OOD (AdversarialQA) F-score of this model of LMIX model addressing *word-dist* bias. Our bias models reach between 18% and 59% of accuracy on easier, i.e., biased data split while between 4% and 19% on the non-biased one.

## B.2.2 Baseline debiasing: Resampling

We train the RESAM analogically to Baseline Fine-tuning experiments (§B.1). Compared to other debiasing methods, RESAM baseline is non-parametric, including no dependence on the bias model.

Even though we find RESAM to be the only method mitigating Prediction bias in all the cases, our further analyses show that its enhancements on OOD datasets vary among biases. Figure 8 shows validation losses from the training on SQuAD re-sampled using RESAM by *word-dist*, while analogically, Figure 9 shows the losses for *sim-ents* bias. While in the former case, RESAM does not stably reach lower loss on OOD datasets, in the latter case, validation losses are consistently lower between steps 7,000 and 8,000, where the SQuAD validation loss used to pick the best-performing model plateaus.

## B.2.3 Learned Mixin

In addition to the implementation and default parameters of Clark et al. (2019a), we find that the additional entropy regularization component  $H$  makes a significant difference in the resulting model evaluation. Therefore we perform a hyperparameter search over the values of  $H$  used for QA by Clark et al. (2019a) on *word-dist* bias, optimizing the OOD performance on AdversarialQA (Bartolo et al., 2020) and eventually fix  $H = 0.4$  over all our experiments.

Following the low initial OOD performance of LMIX as compared to the results of Clark et al. (2019a), we further investigate covariates of this result and identify LMIX’s high sensitivity to bias model; while in the original implementation, TF-IDF similarities of question and answer segment likely never reach 1.0, our generic bias models reaches 1.0 probability for most of the samples marked as biased. Hence, we introduce a parameter of scaling interval  $\langle 0; x \rangle$  of bias model’s scores, where we optimize  $x \in \langle 0.2; 0.4; 0.5; 0.6; 0.7; 0.8; 0.9; 0.95 \rangle$  according to the maximum ID F-score of the debiased model addressing *word-dist* bias, fixing optimal  $x = 0.8$  throughout all other experiments. All other parameters remain identical to the standard fine-tuning (§B.1).

## B.2.4 Confidence Regularization

While the authors of CREG (Utama et al., 2020a) find benefits in its non-parametricity, we find that CREG also shows high sensitivity to a selection of bias model, guiding us to also rescale the prediction of the bias model in the training distillation process. We use the same methodology to pick the scaling interval  $\langle 0; x \rangle$  for CREG as for LMIX and fix  $x =$

1011 0.9 as the optimal one. All other parameters remain  
1012 the identical to the standard fine-tuning (§B.1).

1013 We implement CREG using Transformers library  
1014 (Wolf et al., 2020a) in version 4.19.1.

## 1015 C Exploiting Heuristics Configuration

1016 Here we enumerate the optimal thresholds over  
1017 all pairs of the implemented heuristics, as picked  
1018 according to BERT-BASE-CASED model.

1019 We assess the candidate thresholds among all  
1020 possible values within the range of the computed  
1021 values  $A_h$  computed over  $X = \text{SQuAD}_{\text{valid}}$  (see  
1022 Algorithm 1), with steps of 1 for possible values  
1023 higher than 1 and 0.1 for values between 0 and  
1024 1, within the valid interval; We set the validity  
1025 interval such that the resulting splits of the dataset  
1026 must each have a size of at least two times of the  
1027 sample size parameter, except where there is only  
1028 one significant threshold, and its size is larger than  
1029 the sample size. The optimal threshold value is  
1030 then the one that delivers the highest Prediction  
1031 bias value. We find and use the following optimal  
1032 thresholds of BERT-BASE-CASED evaluated on  
1033  $X = \text{SQuAD}_{\text{valid}}$  for specific biases: 7 for *word-*  
1034 *dist*, 3 for *sim-word*, 4 for *ans-len*, 0.1 for *cos-sim*,  
1035 0 for *sim-ents* and 1 for *subj-pos*. A corresponding  
1036 number of samples in the underperforming groups  
1037 of  $\text{SQuAD}_{\text{valid}}$  ( $n=10,570$ ) are following: 1,651 for  
1038 *word-dist*, 3,281 for *sim-word*, 3,124 for *ans-len*,  
1039 954 for *cos-sim*, 5,006 for *sim-ents* and 1,672 for  
1040 *subj-pos*.

1041 The implementations of some biases’ heuristics  
1042 utilize external libraries for entity recognition or  
1043 TF-IDF vectorization. For these, we used SPACY  
1044 in version 3.4.1 and NLTK in version 3.4.1.

## 1045 D Experimental Environment

1046 Our experiments utilized a single NVidia A100  
1047 GPU with 80 GB of VRAM, a single CPU core,  
1048 and less than 32 GB of RAM. However, all our  
1049 experiments can be run using a lower compute  
1050 configuration, given a longer compute time; The  
1051 inference of a single-sample prediction batch of  
1052 ROBERTA-LARGE as our largest model requires  
1053 only 13 GB of VRAM. The debiasing training  
1054 runs take longer to converge, as compared to stan-  
1055 dard fine-tuning; While the conventional training  
1056 and RESAM converges within 10,000 steps (Fig-  
1057 ures 8 and 9) we find that LMIX requires between  
1058 60,000 and 100,000 steps, and CREG needs be-  
1059 tween 20,000 and 30,000 steps to converge, mak-

1060 ing the debiasing training 4–8 times slower in av-  
1061 erage. In our training configuration, each of the  
1062 reported training runs takes between 50 minutes  
1063 and 1 hour per 10,000 updates. Given that our eval-  
1064 uation already aggregates the bootstrapped results,  
1065 we perform a single run for each experiment, which  
1066 might result in a wider confidence interval and con-  
1067 sistentlly smaller measured volumes of Prediction  
1068 bias.

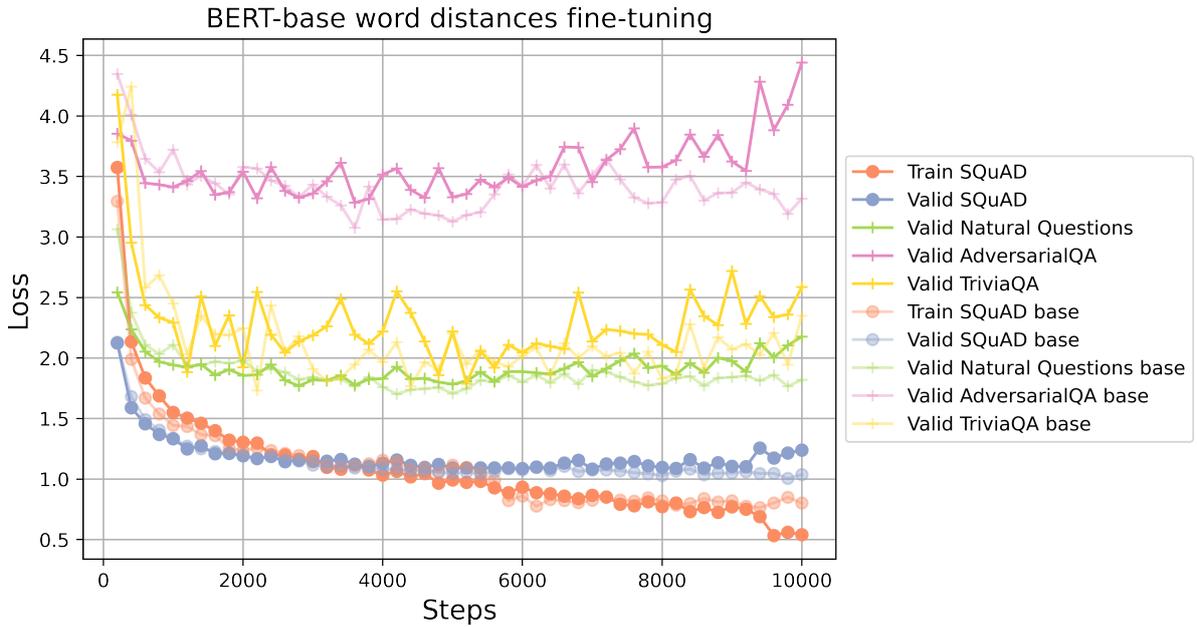


Figure 8: Development of validation loss of **RESAM** addressing *word-dist* bias (darker plots) and standard fine-tuning (lighter plots) for Question Answering on SQuAD, also evaluated on other (OOD) datasets, for the first 10,000 steps.

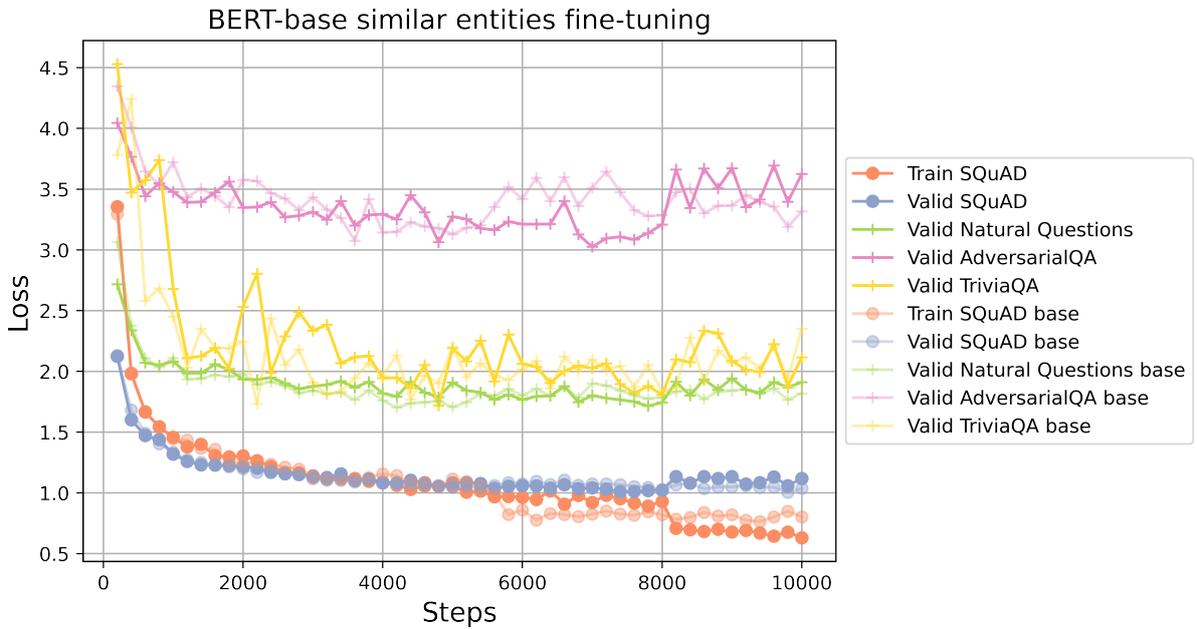


Figure 9: Development of validation loss of **RESAM** addressing *sim-ents* bias (darker plots) and standard fine-tuning (lighter plots) for Question Answering on SQuAD, also evaluated on other (OOD) datasets, for the first 10,000 steps.