

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

Lukáš Mikula^{✉*} and Michal Štefánik^{✉*} and Marek Petrovič^{✉*} and Petr Sojka^{✉*}

[✉]Faculty of Informatics,
Masaryk University, Czech Republic

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 while existing debiasing methods *can* mitigate reliance on a chosen spurious feature, the OOD performance gains of these methods can *not* be explained by mitigated reliance on biased features, suggesting that biases are *shared* among different QA datasets. Finally, we evidence this to be the case by measuring that performance of models trained on different QA datasets rely on bias features *comparably* to the ID model. We hope these results will motivate future work to refine the reports of LMs' robustness to a level of adversarial samples addressing specific spurious features.

1 Introduction

Unsupervised pre-training and vast parametrization (Devlin et al., 2019; 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.,

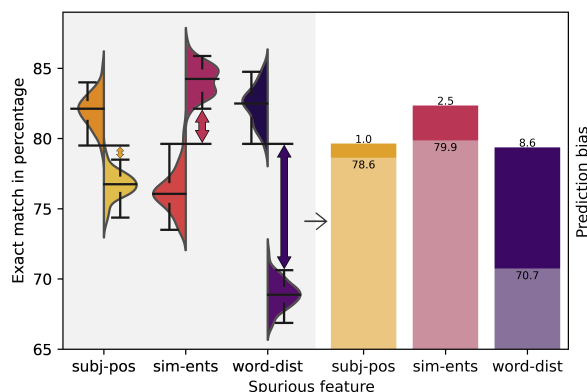


Figure 1: We quantify model reliance on a spurious feature using bootstrapped evaluation on segments of data separated by exploiting chosen bias (left) and subsequently, by measuring the difference in model's performance over these two groups (right), that we refer to as *Prediction bias* (§3).

2020) in Natural Language Inference (NLI) or the intersection of the question and answer vocabulary (Shinoda et al., 2021) in extractive Question Answering (QA).

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 (Mehrabani 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

*First two authors contributed equally

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

We find that avoiding reliance on spurious features does *not imply* improvements in OOD performance; in many cases, debiasing methods mitigate the model’s prediction bias, but the OOD performance drops, while counterintuitively, a magnification of bias reliance can also bring large OOD gains. Aiming to explain this, we directly evaluate the prediction bias of models trained on different datasets and confirm that even models trained on OOD datasets often rely on the *same* spurious correlations *comparably* to the ID models. This finding motivates the presented assessment of model robustness towards known biases, in addition to OOD performance.

This paper is structured as follows. Section 2 overviews data biases observed in NLP datasets, recent debiasing methods, and the previous methods related to measuring inclination to spurious correlations. Section 3 presents our method for measuring the significance of specific biases. We follow in Section 4 with details on our evaluation setup, including the tested debiasing methods, addressed bias features, and the design of a set of heuristics that can exploit them. Subsequently, in Section 5, we measure and report models’ robustness to biases and OOD datasets before and after applying the selected debiasing methods and wrap up our observations in Sections 6 and 7.

Problem definition Given a set of inputs $X = x_{1..i}$ with corresponding labels $Y = y_{1..i}$ from a dataset \mathcal{D}_{ID} , a model M learns a *task* \mathcal{T} by identifying *features* $\mathcal{F}_{1..n}$ that map each x_j to a corresponding y_j , assuming that the learned features must be *consistent* with \mathcal{D}_{ID} . Since the learned $\mathcal{F}_{1..n}$ are distributed in M and can not be directly evaluated, we assess whether the learned features are *robust* for the task \mathcal{T} by evaluating M on samples X_{OOD} of the same task, but drawn from $\mathcal{D}_{OOD} \not\approx \mathcal{D}_{ID}$; we assume that if $\mathcal{F}_{1..n} \in M$ are robust, the model will also perform well on X_{OOD} . However, the consistency of the learned \mathcal{F}_k with both X_{ID} and X_{OOD} is merely a *necessary* and not a *sufficient* condition

for \mathcal{F}_k to be robust; If there exists a pair (x, y) such that the pair is a *valid* sample of the task \mathcal{T} , but is not consistent with \mathcal{F}_k , we denote \mathcal{F}_k as *spurious* or *bias features* for \mathcal{T} and refer to models’ reliance on such features as *prediction bias*.

2 Background

Spurious correlations of NLP datasets Previous work analyzing LLMs’ error cases identified numerous false assumptions that LLMs use in prediction and can be misused to notoriously draw wrong predictions with the model.

In Natural Language Inference (NLI), where the task is to decide whether a pair of sentences entail one another, McCoy et al. (2019) identify LLMs’ reliance on a lexical overlap and on specific shared syntactic units such as the constituents in the processed sentence pair. Asael et al. (2022) identify the model’s sensitivity to meaning-invariant structure permutations. Similarly, Chaves and Richter (2021) identify BERT’s reliance on the invariant morpho-syntactic composition of the input.

In Question Answering, LLMs often rely on the positional relation of the question and possible answer words, such as assuming their close *proximity* (Jia and Liang, 2017). Bartolo et al. (2020) find that models tend to assume that questions and answers contain similar *keywords*, remaining vulnerable to samples with none or multiple occurrences of the keywords in the context. Ko et al. (2020) show models’ preference for the answers in the first two sentences of the context, being statistically most likely to answer human-curated questions.

A perspective direction circumventing the biases introduced in data collection is presented in adversarial data collection (Jia and Liang, 2017; Bartolo et al., 2020) where the annotators collect the dataset with the intention of fooling the likely-biased model, possibly enhancing the model-in-the-loop in several fine-tuning iterations. Still, some doubts remain, as other work provides evidence that models trained on adversarial data may work better on adversarial datasets but underperform on other datasets (Kaushik et al., 2021), or introduce its own set of biases (Kovatchev et al., 2022). Nevertheless, our experiments (§5.2) show that training models on an adversarially-collected AdversarialQA dataset turns out to be among the most effective approaches to mitigating known prediction biases in question answering.

Debiasing methods A well-established line of work proposes to address the known dataset biases in the training process. Karimi Mahabadi et al. (2020) and He et al. (2019) obtain a more robust, debiased model by (i) training a *biased model* that exploits the unwanted bias, followed by (ii) training the debiased model as a complement to the biased one in a Product-of-Experts (PoE) framework (Hinton, 2002). Clark et al. (2019a) extend this framework in the LearnedMixIn method, learning to weigh the contribution of the biased and debiased model in the complementary ensemble. Niu and Zhang (2021) simulate the model for non-biased, out-of-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 bias-detection model might bring further benefits to this framework, consistently with our observations (§6). 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, 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

```

func measure_bias(M, X, h, Th):
    Ah ← h(X)
    X1 ← x1 ∈ X : Ah(x1) ≤ Th
    X2 ← x2 ∈ X : Ah(x2) > Th
    foreach X1′ ∈ repeat(sample(X1)) do
        | E1 ← E1 + evaluate(M(X1′))
    foreach X2′ ∈ repeat(sample(X2)) do
        | E2 ← E2 + evaluate(M(X2′))
    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*₁ and *X*₂) obtained by segmenting the samples of *X* by the *attribute* *A_h* = *h*(*X*) on a given threshold *T_h*.

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

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. Further, compared to the previous work, we assess models’ reliance on a range of 7 spurious features, making our overall conclusions more robust.

An ability to measure a model’s reliance on undesired features is also applicable in quantifying socially problematic biases. Previous work also utilizes specialized domain knowledge in models’ 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 LMs’ 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 visualized in Figure 1, described in

Algorithm 1 and can be used to measure biases of any other QA model using the project repository¹.

We start by (i) implementing a *heuristic*, i.e. a method $h : X \rightarrow \mathbb{R}$, that for *all* samples of dataset X computes an *attribute* $A_h \in \mathbb{R}$ corresponding to the feature \mathcal{F} that we suspect as non-representative, yet predictive for our training set and (ii) we compute $h(x)$ for each sample x of evaluation dataset X . (iii) We choose a threshold T_h that we use to (iv) split the dataset into two segments by A_h . Finally, (v) we evaluate the assessed model M on *both* of these segments, in our case using Exact match evaluation, 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 the 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 down 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 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 can be present in datasets. In Question answering, previous work identifies several spurious covariates in the SQuAD dataset (Rajpurkar et al., 2016); we build upon these findings and further extend the list of covariates learnable from SQuAD.

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 +.

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

¹<https://github.com/MIR-MU/isbiased>

sentences. The exploiting heuristic first segments the context into sentences, and 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: To (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

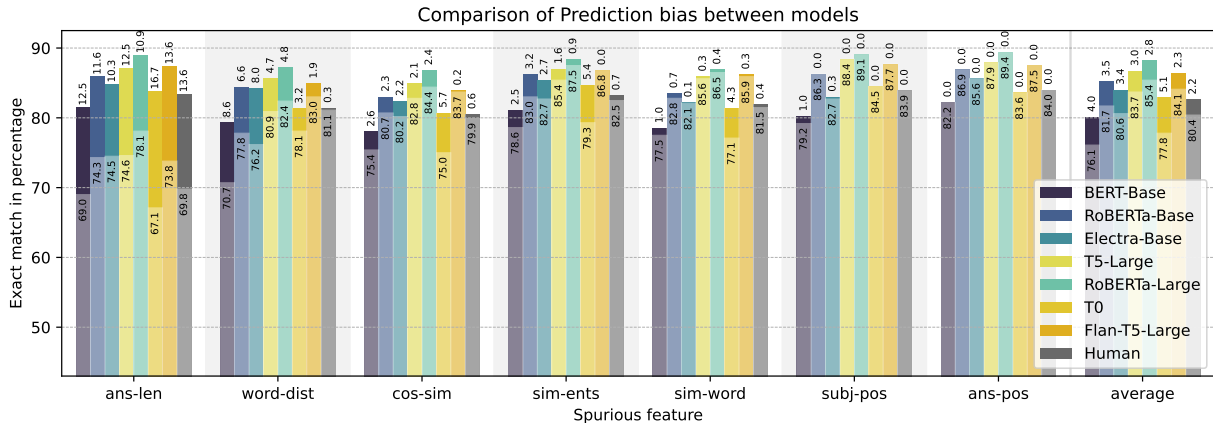


Figure 2: **Prediction bias per pre-trained model.** The worse-performing split performance (lower bars) and Prediction bias (upper bars, sorted by group average) of QA models trained from different pre-trained LLMs, trained and evaluated on SQuAD for Exact match. Per-group bootstrapping of 100 repeats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

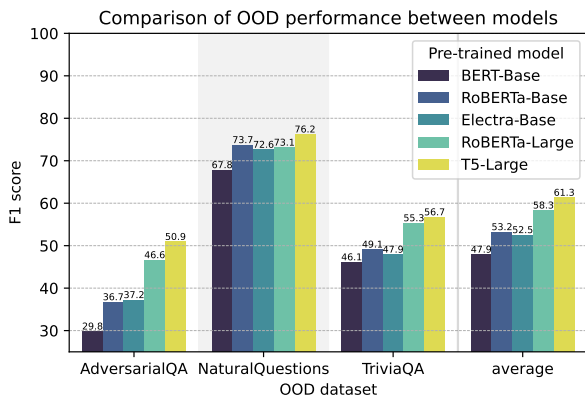


Figure 3: **OOD performance per pre-trained model.** Comparison of F1-score of different models fine-tuned on SQuAD and evaluated on listed OOD datasets.

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 custom-tailored 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 choose 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).

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 a 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.

Analogically, 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 the exception of generative fine-tuning (T5), which outperforms others on OOD datasets but not on a reduction of the reliance on spurious features.

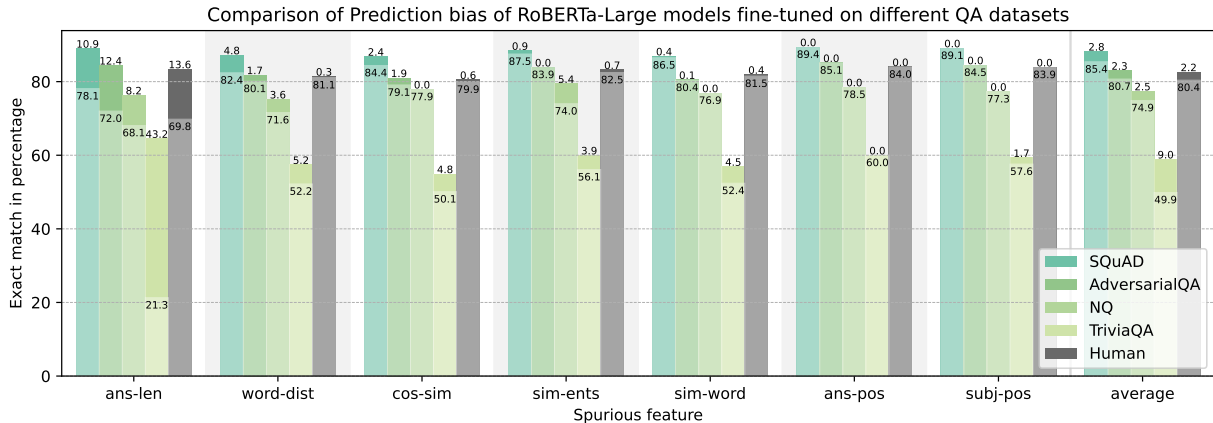


Figure 4: **Prediction bias per dataset.** The worse-performing split performance (lower bars) and Prediction bias (upper bars) of RoBERTa-LARGE trained on different QA datasets, evaluated on a validation split of SQuAD for Exact match. All evaluation splits are identical, identified as maximal for the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

5.2 Prediction bias of OOD models

Figure 4 compares Prediction bias over RoBERTa-LARGE models trained on different datasets. All evaluations are split on heuristics’ thresholds T_h optimal for the 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 learned on SQuAD are also learned from at least one OOD dataset. For the Trivia model, *all* types of biases identified in SQuAD are magnified.

We specifically note the comparison of the Prediction bias of the SQuAD model to the model trained on AdversarialQA, collected adversarially to a SQuAD model. We find that the 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*, 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: **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 in **bold**.

	Original model			29.8 / 67.8 / 46.1					
	ReSam			LMix			CReg		
	<i>AQA</i>	<i>NQ</i>	<i>Trivia</i>	<i>AQA</i>	<i>NQ</i>	<i>Trivia</i>	<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	+1.4 / +7.5 / -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		

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 average F1-score on OOD datasets by 2.8% and by 7.5 specifically on *NaturalQuestions*, but the Prediction bias of such model increases by 1.1 points. Similarly, CREG delivers 1.5-point average gain of F1-score on OOD when addressing *sim-word* bias but this also 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, increas-

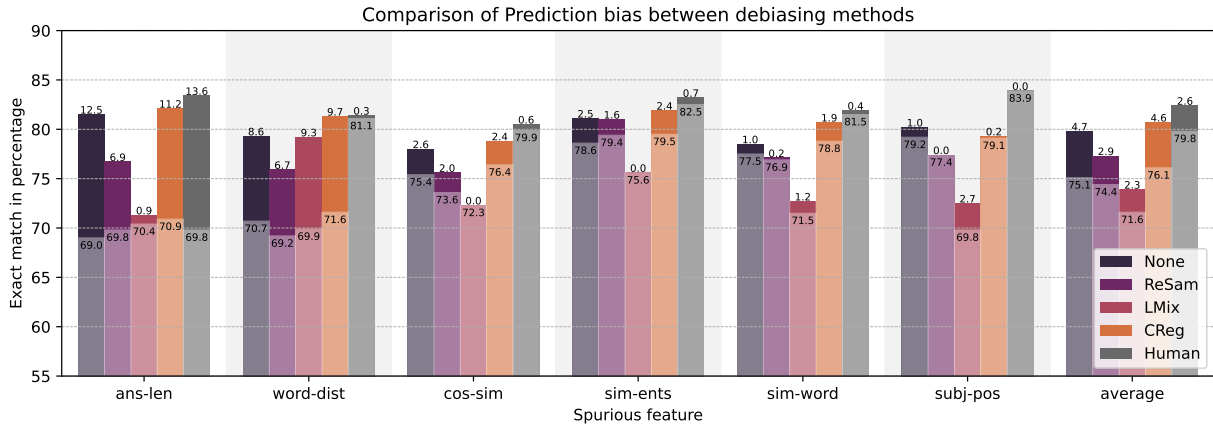


Figure 5: **Prediction bias per debiasing methods.** The worse-performing split performance (lower bars) and Prediction bias (upper bars) of BERT-BASE trained using selected debiasing methods, evaluated for Exact match on validation SQuAD. Per-group evaluations were measured using bootstrapping of 100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

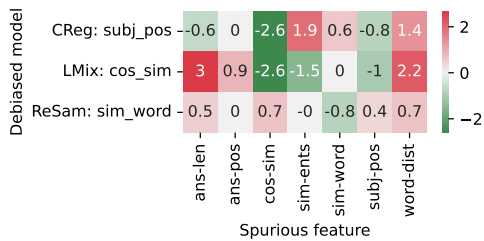


Figure 6: **Cross-bias evaluation of debiased models.** A relative change of Prediction bias by all spurious correlations, caused by applying inspected debiasing methods on BERT-BASE QA model, in addressing specified spurious correlation. A full matrix is in Appx. A, Fig. 7.

ing 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, To’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 biased features. (ii) In practice, we find cases 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 the existing and to identify 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* minimizing 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 learnt from SQuAD.

We hope that our analyses will motivate future work to assess models' robustness also on a more detailed level of specific bias features, evading false conclusions on models' robustness, and, ultimately, accelerating progress towards creating more robust and reliable language models.

Limitations

We highlight the limitation of our proposed evaluation method in the non-trivial *interpretation* of the measured results, which we discuss in Section 3; We propose to measure the models' reliance on a bias feature as a difference of *confidence intervals* of model performance on two data splits. This makes the conclusions about models' reliance (vs non-reliance) on a biased feature more robust, but

it also perplexes the interpretation of measured absolute values. As a consequence, in the cases of different bias features ($\mathcal{F}_1, \mathcal{F}_2$) with very close prediction bias values, one should restrain from statements such as "model M is more biased towards \mathcal{F}_1 than \mathcal{F}_2 ".

We also underline that some biased features correlate with a *natural* difference in the samples' difficulty. In such settings, a polarization of model performance might not be caused by its reliance on the spurious feature, but rather by other, natural features of the task. To disentangle the model's over-reliance on a biased feature from other aspects, we recommend contextualizing measured prediction bias with additionally measuring a *human* level of prediction bias, that can be assessed on a set of duplicate annotations.

In our experiments, we measured considerable differences in natural difficulty only for a single feature – answer length – where it is likely more difficult to delimit the answer span for longer answers properly. We find that most models rely on this feature comparably to humans and refine our conclusions in Section 5.1 accordingly.

Acknowledgments

We acknowledge the Centre for Biomedical Image Analysis at Masaryk University supported by MEYS CR (LM2023050 and CZ.02.1.01/0.0/0.0/18_046/0016045 Czech-BioImaging) for their support in creating the models evaluated within this paper.

References

- Dimion Asael, Zachary Ziegler, and Yonatan Belinkov. 2022. [A Generative Approach for Mitigating Structural Biases in Natural Language Inference](#). In *Proceedings of the 11th Joint Conference on Lexical and Computational Semantics*, pages 186–199, Seattle, Washington. ACL.
- Max Bartolo, A Roberts, Johannes Welbl, Sebastian Riedel, and Pontus Stenetorp. 2020. [Beat the ai: Investigating adversarial human annotation for reading comprehension](#). *Transactions of the Association for Computational Linguistics*, 8:662–678.
- Shikha Bordia and Samuel R. Bowman. 2019. [Identifying and Reducing Gender Bias in Word-Level Language Models](#). In *Proceedings of the 2019 Conference of the North American Chapter of the ACL: Student Research Workshop*, pages 7–15, Minneapolis, Minnesota. ACL.

- Rui P. Chaves and Stephanie N. Richter. 2021. [Look at that! BERT can be easily distracted from paying attention to morphosyntax.](#) In *Proceedings of the Society for Computation in Linguistics 2021*, pages 28–38. ACL.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. [Scaling Instruction-Finetuned Language Models.](#) *arXiv e-prints*, page arXiv:2210.11416.
- Christopher Clark, Mark Yatskar, and Luke Zettlemoyer. 2019a. [Don't take the easy way out: Ensemble based methods for avoiding known dataset biases.](#) In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4069–4082, Hong Kong, China. Association for Computational Linguistics.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019b. [What does BERT look at? an analysis of BERT's attention.](#) In *Proc. of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. ACL.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. [ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators.](#) *CoRR*, abs/2003.10555v1.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.](#) In *Proc. of the 2019 Conference of the NAACL: Human Language Technologies*, pages 4171–4186, Minneapolis, USA. ACL.
- 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. ACL.
- He He, Sheng Zha, and Haohan Wang. 2019. [Unlearn Dataset Bias in Natural Language Inference by Fitting the Residual.](#) In *Proc. of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019)*, pages 132–142, Hong Kong, China. ACL.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. [Distilling the Knowledge in a Neural Network.](#) NIPS 2014 Deep Learning Workshop.
- Geoffrey E. Hinton. 2002. [Training Products of Experts by Minimizing Contrastive Divergence.](#) *Neural Computation*, 14(8):1771–1800.
- Matthew Honnibal and Ines Montani. 2017. [spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing.](#)
- Robin Jia and Percy Liang. 2017. [Adversarial examples for evaluating reading comprehension systems.](#) In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2021–2031, Copenhagen, Denmark. Association for Computational Linguistics.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. [Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension.](#) *arXiv preprint arXiv:1705.03551*.
- 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 ACL*, pages 8706–8716. ACL.
- Divyansh Kaushik, Douwe Kiela, Zachary C. Lipton, and Wen-tau Yih. 2021. [On the Efficacy of Adversarial Data Collection for Question Answering: Results from a Large-Scale Randomized Study.](#) In *Proc. of the 59th Annual Meeting of the ACL and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6618–6633. ACL.
- Miyoung Ko, Jinhyuk Lee, Hyunjae Kim, Gangwoo Kim, and Jaewoo Kang. 2020. [Look at the first sentence: Position bias in question answering.](#) *arXiv preprint arXiv:2004.14602*.
- Venelin Kovatchev, Trina Chatterjee, Venkata S Govindarajan, Jifan Chen, Eunsol Choi, Gabriella Chronis, Anubrata Das, Katrin Erk, Matthew Lease, Junyi Jessy Li, Yating Wu, and Kyle Mahowald. 2022. [longhorns at DADC 2022: How many linguists does it take to fool a Question Answering model? A systematic approach to adversarial attacks.](#) In *Proceedings of the First Workshop on Dynamic Adversarial Data Collection*, pages 41–52, Seattle, WA. ACL.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. 2019. [Natural questions: a benchmark for question answering research.](#) *Transactions of the Association for Computational Linguistics*, 7:453–466.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, M. Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [RoBERTa: A Robustly Optimized BERT Pretraining Approach.](#) *CoRR*.

- Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. [Right for the Wrong Reasons: Diagnosing Syntactic Heuristics in Natural Language Inference](#). In *Proc. of the 57th Annual Meeting of the ACL*, pages 3428–3448, Florence, Italy. ACL.
- Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. [A Survey on Bias and Fairness in Machine Learning](#). *ACM Comput. Surv.*, 54(6).
- Yulei Niu, Kaihua Tang, Hanwang Zhang, Zhiwu Lu, Xian-Sheng Hua, and Ji-Rong Wen. 2021. [Counterfactual VQA: A Cause-Effect Look at Language Bias](#). In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 12700–12710. Computer Vision Foundation / IEEE.
- Yulei Niu and Hanwang Zhang. 2021. [Introspective distillation for robust question answering](#). In *Advances in Neural Information Processing Systems*, volume 34, pages 16292–16304. Curran Associates, Inc.
- Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. 2022. [BBQ: A hand-built bias benchmark for question answering](#). In *Findings of the ACL: ACL 2022*, pages 2086–2105, Dublin, Ireland. ACL.
- Alec Radford and Karthik Narasimhan. 2018. [Improving Language Understanding by Generative Pre-Training](#).
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(146):1–67.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ Questions for Machine Comprehension of Text](#). In *Proc. of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, USA. ACL.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. 2020. [Distributionally Robust Neural Networks](#). In *International Conference on Learning Representations*.
- Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. 2022. [Multitask Prompted Training Enables Zero-Shot Task Generalization](#). In *International Conference on Learning Representations*.
- Kazutoshi Shinoda, Saku Sugawara, and Akiko Aizawa. 2021. [Can Question Generation Debias Question Answering Models? A Case Study on Question–Context Lexical Overlap](#). In *Proceedings of the 3rd Workshop on Machine Reading for Question Answering*, pages 63–72, Punta Cana, Dominican Republic. ACL.
- Michal Štefánik, Vít Novotný, Nikola Groverová, and Petr Sojka. 2022. [Adaptor: Objective-Centric Adaptation Framework for Language Models](#). In *Proceedings of the 60th Annual Meeting of the ACL: System Demonstrations*, pages 261–269, Dublin, Ireland. ACL.
- Lifu Tu, Garima Lalwani, Spandana Gella, and He He. 2020. [An Empirical Study on Robustness to Spurious Correlations using Pre-trained Language Models](#). *Transactions of the ACL*, 8:621–633.
- Prasetya Ajie Utama, Nafise Sadat Moosavi, and Iryna Gurevych. 2020a. [Mind the trade-off: Debiasing NLU models without degrading the in-distribution performance](#). In *Proc. of the 58th Annual Meeting of the ACL*, pages 8717–8729. ACL.
- Prasetya Ajie Utama, Nafise Sadat Moosavi, and Iryna Gurevych. 2020b. [Towards Debiasing NLU Models from Unknown Biases](#). In *Proc. of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7597–7610. ACL.
- Thomas Wolf, Julien Chaumond, Lysandre Debut, Victor Sanh, Clement Delangue, Anthony Moi, Pierric Cistac, Morgan Funtowicz, Joe Davison, Sam Shleifer, et al. 2020a. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020b. [Transformers: State-of-the-Art Natural Language Processing](#). In *Proc. of the 2020 Conf. EMNLP: System Demonstrations*, pages 38–45. ACL.
- Mingzhu Wu, Nafise Sadat Moosavi, Andreas Rücklé, and Iryna Gurevych. 2020. [Improving QA Generalization by Concurrent Modeling of Multiple Biases](#). In *Findings of the ACL: EMNLP 2020*, pages 839–853. ACL.
- Ruibin Xiong, Yimeng Chen, Liang Pang, Xueqi Cheng, Zhi-Ming Ma, and Yanyan Lan. 2021. [Uncertainty calibration for ensemble-based debiasing methods](#). In *Advances in Neural Information Processing Systems*.

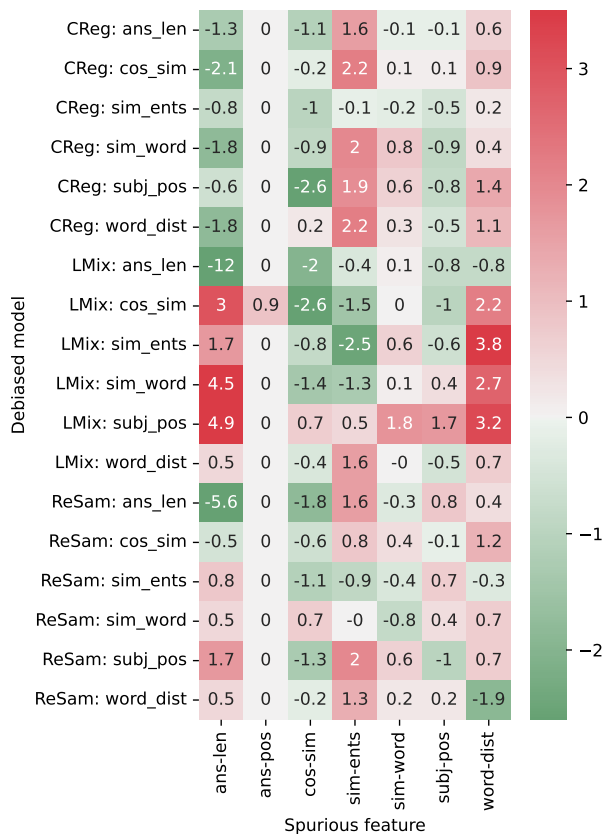


Figure 7: **Full cross-bias evaluation of debiased models.** A relative change of Prediction bias by all spurious correlations, caused by applying inspected debiasing methods on BERT-BASE QA model, in addressing specified spurious correlation.

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

Chunting Zhou, Xuezhe Ma, Paul Michel, and Graham Neubig. 2021. **Examining and Combating Spurious Features under Distribution Shift.** In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 12857–12867. PMLR.

A Cross-Bias Matrix of All Debiased Models

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

B Details of Training Configurations

This section overviews all configurations that we have set in training the debiased models (§4.3 – 4.4)

as well as the conventional QA fine-tuning comparing the impact of pre-training on QA models’ robustness (§4.2).

B.1 Standard Fine-tuning

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

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

B.2 Debiasing Training Experiments

B.2.1 Bias models

The canonical debiasing implementations utilize bias-specific models for identifying bias; Clark 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 biased 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 perfor-

mance 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 resampled 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 $x = 0.9$ as the optimal one. All other parameters remain identical to the standard fine-tuning (§B.1).

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

C Exploiting Heuristics Configuration

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

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

954 for *cos-sim*, 5,006 for *sim-ents* and 1,672 for *subj-pos*.

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

D Experimental Environment

Our experiments utilized a single NVidia A100 GPU with 80 GB of VRAM, a single CPU core, and less than 32 GB of RAM. However, all our experiments can be run using a lower compute configuration, given a longer compute time; The inference of a single-sample prediction batch of RoBERTA-LARGE as our largest model requires only 13 GB of VRAM. The debiasing training runs take longer to converge, as compared to standard fine-tuning; While the conventional training and RESAM converge within 10,000 steps (Figures 8 and 9) we find that LMIX requires between 60,000 and 100,000 steps, and CREG needs between 20,000 and 30,000 steps to converge, making the debiasing training 4–8 times slower in average. In our training configuration, each of the reported training runs takes between 50 minutes and 1 hour per 10,000 updates. Given that our evaluation already aggregates the bootstrapped results, we perform a single run for each experiment, which might result in a wider confidence interval and consistently smaller measured volumes of Prediction bias.

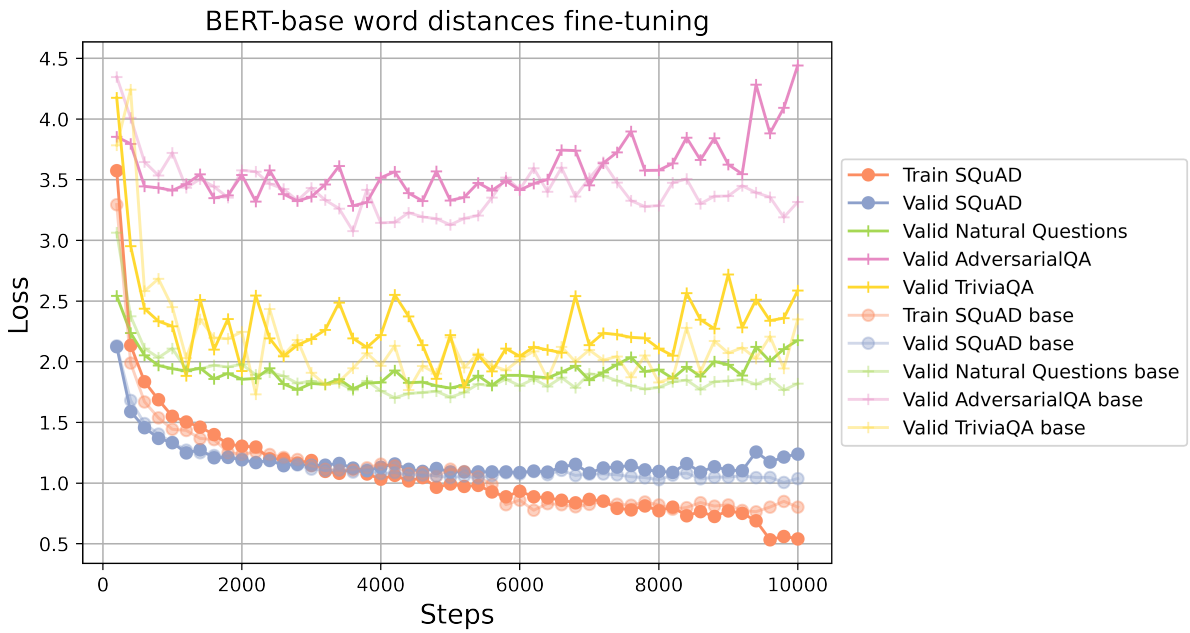


Figure 8: Development of validation loss of **RESAM** addressing *word-dist* bias (darker plots) and standard fine-tuning (lighter plots) for Question Answering on SQuAD, also evaluated on other (OOD) datasets, for the first 10,000 steps.

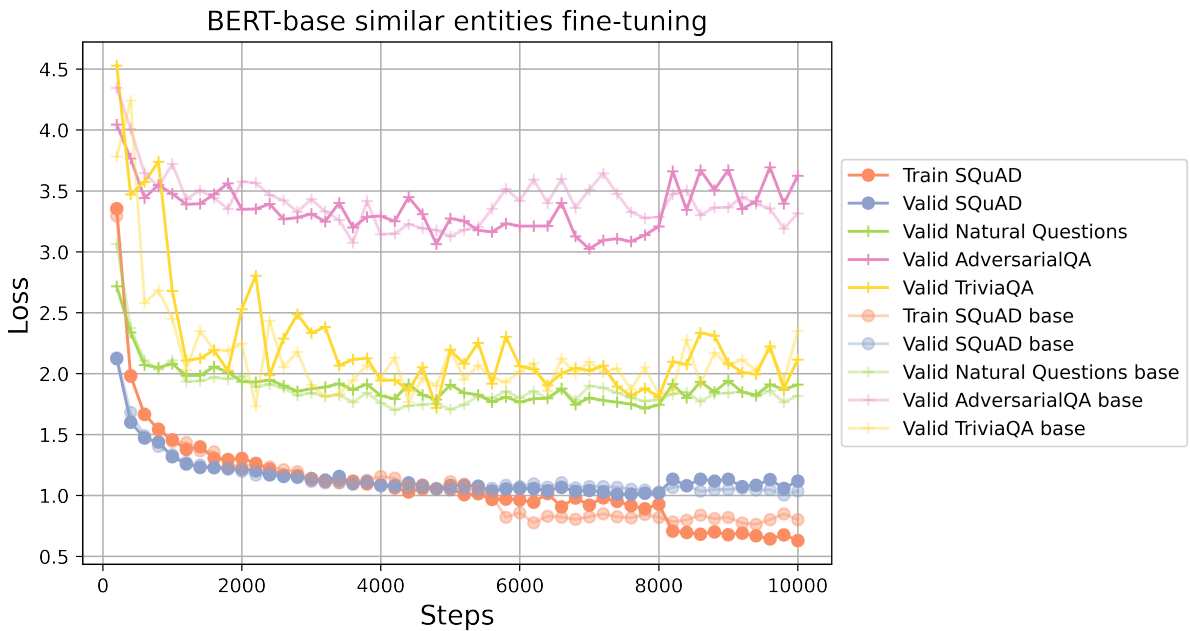


Figure 9: Development of validation loss of **RESAM** addressing *sim-ents* bias (darker plots) and standard fine-tuning (lighter plots) for Question Answering on SQuAD, also evaluated on other (OOD) datasets, for the first 10,000 steps.