

An Empirical Study on Robustness of Language Models via Decoupling Representation and Classifier

Anonymous ACL submission

Abstract

Recent studies indicate that shortcut learning behavior exists in language models, and thus a number of mitigation methods are proposed, such as advanced PLMs and debiasing methods. However, few studies have explored how different factors affect the robustness of language models. To bridge this gap, we study the different PLMs and analyze the effect of representations and classifiers on robustness using probing techniques on the NLU tasks. First, we find that the low robustness of language models is not due to the inseparability of representations on the challenging dataset. Second, we find that a potential reason for the difficulty in improving the robustness of language models is the significantly high similarity between the representations with opposite semantics from in-distribution and out-of-distribution. Third, we find that debiasing methods are likely to distort representations and merely improve performance by better classifiers in some cases¹. Finally, we propose a probing tool to measure the impact on the robustness of language models from representations and classifiers using the decoupled training strategy with debiasing methods. In addition, we conduct extensive experiments on real-world datasets, suggesting the effectiveness of the proposed methods.

1 Introduction

Pre-trained Language Models (PLMs), such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), have achieved state-of-the-art results for Natural Language Understanding (NLU) tasks. Despite their successes, recent studies show the phenomenon that PLMs are prone to learning superficial surface patterns that are spuriously associated with the target label, and to make use of biases/artifacts from the dataset as shortcuts for prediction (Gururangan et al., 2018; McCoy et al.,

¹In this work, we denote the **representation** as the output of PLMs and **classifier** as the several fully connected layers that are used to serve the classification purpose.

2019; Utama et al., 2020b), which is defined as shortcut learning (Geirhos et al., 2020). For a common NLU task, Natural Language Inference (NLI), the shortcut learning behavior is defined as that the model achieves high accuracy only by using specific words but not understanding the language (Naik et al., 2018; Sanchez et al., 2018; Du et al., 2021). As a result, the models perform poorly on out-of-distribution (OOD) examples.

The quality of representations is widely considered to be the key reason for the shortcut learning and poor generalization ability of the NLU models. Because of this, there is a large body of literature that analyzes and understands the learned representation (Perone et al., 2018; Krasnowska-Kieraś and Wróblewska, 2019; Pruksachatkun et al., 2020; Mendelson and Belinkov, 2021). Unlike previous work, we intend to answer the following research questions: 1) *Whether the low robustness of language models is primarily due to representations not being easily separable?* 2) *What is the relative role of the representation and the classifier in the low robustness of language models?*

In this work, we apply the t-SNE visualization technique and DIRECTPROBE (Zhou and Srikumar, 2021) probing technique to representations that are from five PLMs: BERT, RoBERTa, BART (Lewis et al., 2020), ELECTRA (Clark et al., 2020b), and DeBERTa (He et al., 2021). Meanwhile, we present a new probing strategy based on debiasing methods. Based on the above techniques and strategies, our findings on the robustness of language models are briefly described below.

To begin with, we visualize the representations of the above five pre-trained language models. We find that not only [CLS] but also [MEAN] embeddings form clearer boundaries between representations of different labels, despite only [CLS] embeddings are fed into the classifiers (§4.1).

Then, we investigate the geometric structure of representations. Via DIRECTPROBE probing tech-

nique, we obtain clusters with only the same label representations in each cluster. We find that the number of clusters is equal to the number of label categories in most cases, which means that the representations are linearly separable (§4.2).

Furthermore, we further study the properties of representations by computing the similarity between the clusters with opposite semantics on MNLI and HANS. We find that a possible reason why it is not easy to improve the robustness of models is that the representations with opposite semantic labels are too similar (§4.3).

Finally, we investigate the effect of encoders and classifiers on the robustness of language models respectively. Using debiasing methods as a probing tool, we find that both the representation and the classifier of the models play a significant role in the shortcut learning behavior for the NLI task. Furthermore, we find that debiasing methods do not always improve the quality of representations. Instead, they only improve performance by optimizing the classifiers in some cases (§4.4).

2 Preliminaries

In this work, we probe the representations and classifiers of PLMs after fine-tuning. To this end, we briefly introduce two techniques that we use in our analyses: the probing method and the ensemble-based debiasing framework.

2.1 Probing Method

Normally, trained classifiers are used as probes to understand the quality of the information encoded in the representation, which are trained with the encoders frozen (Hewitt et al., 2021; Whitney et al., 2021; Belinkov, 2021). However, the classifier probes focus only on the performance of the target task and cannot clarify the representation in detail (Zhou and Srikumar, 2022). To this end, we apply a probing technique named DIRECTPROBE (Zhou and Srikumar, 2021) instead of classifier probes. It can provide a fine-grained analysis of the representation from a geometric perspective.

The representation in the form of embeddings is fed into DIRECTPROBE, and the number of clusters is returned. These clusters satisfy the condition that the example points contained in a cluster must have the same label and that there are no overlaps between any two clusters (that is, there exists a separator between the two sets of example points). We can learn about various linguistic attributes of

the representation by measuring the properties of the corresponding clusters. In this work, we focus on an important property: the number of clusters.

Number of Clusters The number of clusters can quantify the linear separability of the representation for a special task. In particular, when the number of clusters equals the number of label categories, the embeddings of examples with the same label are close enough in the semantic representation space. In this ideal case, the models can achieve perfect performance with a simple linear classifier. In contrast, when the number of clusters is more than the number of label categories, the example points with the same label are grouped into at least two clusters. This suggests a complex geometric structure of representations, and a complex classifier is needed to achieve desirable performance.

2.2 Ensemble-based Debiasing Framework

The ensemble-based debiasing framework (EBD) (Xiong et al., 2021) is generally used to mitigate the shortcut learning behavior of the NLU models. This framework has the advantage that it is a model-agnostic debiasing framework, which makes it possible to debias models adaptively. EBD framework consists of bias-only models and debiasing methods. Bias-only models are used to make the main models perform debiasing training by adjusting the learning target. Debiasing methods provide strategies on how to debias the main models in practice. The EBD framework is commonly formalized as a two-stage method (Clark et al., 2019; Sanh et al., 2021). In the first stage, the bias-only model is trained to recognize simple and hard examples. In the second stage, the main models are trained as an ensemble with the bias-only model according to the selected debiasing methods.

2.2.1 Bias-only Model

Recently, several works in the literature have proposed exploring bias-only models to improve the performance of the EBD framework. For example, Utama et al., 2020b, Sanh et al., 2021, and Clark et al., 2020a try to reduce the need for a prior knowledge on bias or shortcut. In these works, they obtain bias-only models with two different strategies: i) training a copy of the main model with a small random subset of training examples for a few epochs; and ii) using a shallow or small model with limited capacity. In our work, we take the first strategy, and more details are given in Appendix A.

In the following, we describe the workflow of bias-only models. For clarity, we denote the bias-only model by f_b . Given an example (x^i, y^i) in the training dataset, we denote the output of f_b as $f_b(x^i) = p_b^i$. Probability p_b^i can quantify how much the model learns about shortcut features from example (x^i, y^i) (i.e., how likely this example contains biases). Specifically, the extent to which models learn shortcut features can be evaluated by $p_b^{(i,c)}$ which denotes the probability of p_b^i on the label y^i , where c is the index of the correct category in the label y^i . For example, when $p_b^{(i,c)}$ is closer to 1 (i.e., the bias-only model is more confident about the example x^i on the label y^i), the model learns more potentially shortcut features. Instead, when $p_b^{(i,c)}$ is closer to 0, the bias-only model is more unconfident about the example x^i on the label y^i . As such, the example x^i is likely to be a hard example to which the model is supposed to pay more attention during training.

2.2.2 Debiasing Method

We first denote the main model by f_d parameterized by θ_d , and then use the bias-only model f_b obtained in §2.2.1 to perform debiasing training on f_d . In this work, we mainly investigate two common model-agnostic debiasing methods: sample re-weighting (Schuster et al., 2019) and product-of-experts (Clark et al., 2019; Karimi Mahabadi et al., 2020). In the following, we describe the implementation details of these two methods.

Example Re-weighting Example re-weighting is a simple yet effective debiasing method. It can be briefly summarized as re-weighting the importance of a given training example (x^i, y^i) by directly assigning a weight to the example (x^i, y^i) . The weight is formalized as $1 - p_b^{(i,c)}$. Thus, the individual loss of the example (x^i, y^i) for the parameters θ_d is defined as follows:

$$\mathcal{L}(\theta_d) = -\left(1 - p_b^{(i,c)}\right) y^i \cdot \log p_d,$$

where p_d is the *softmax* output of the main model f_d . Here, we regard training samples with high probability by the bias-only model as biased/shortcut samples. When the bias-only model assigns a high probability to $p_b^{(i,c)}$, the contribution of a training example to $\mathcal{L}(\theta_d)$ is reduced.

Product-of-Experts In this method, the main model (i.e., a debiased model) is trained in an ensemble with a bias-only model. Specifically, the

softmax outputs of the main model f_d and the bias-only model f_b are combined to form new predictions. Then they are used to calculate the new loss while optimizing the parameters θ_d . The individual loss of the example (x^i, y^i) for the parameters θ_d is defined as follows:

$$\mathcal{L}(\theta_d) = -y^i \cdot \log \text{softmax}(\log p_d + \log p_b).$$

During training with debiasing methods, the parameters of the bias-only model f_b are frozen to lower the importance of biased examples in training loss, and only the parameters of the main model f_d are optimized. During the inference time, only the prediction probability of f_d is used.

3 Experimental Setup

3.1 Tasks & Datasets

In this work, we focus on a common NLU task: Natural Language Inference (NLI), the model of which is presented with a pair of sentences and asked to return the relationship between their meanings (Williams et al., 2018). A pair of sentences contains a premise sentence and a hypothesis sentence. The relationship between their meanings is one label of *entailment*, *neutral*, and *contradiction*.

MNLI MNLI (Williams et al., 2018) is divided into training dataset, matched development dataset, and mismatched development set. The training dataset and the matched development dataset are derived from the same five genres, and the mismatched development dataset are derived from the other five genres. Typically, we first use the training dataset to train an NLU model, which consists of 392,702 instances. Then, we use the matched development dataset to choose an optimal NLU model, which consists of 9,815 instances.

HANS HANS (McCoy et al., 2019) has been proposed to evaluate whether models have learned statistical patterns or semantic understanding and reasoning. It focuses on three heuristics: the lexical overlap heuristic, the subsequence heuristic, and the constituent heuristic. HANS consists of 30,000 synthetic instances, and distributes 10,000 ones to each of the heuristics. Note that HANS is only used to evaluate the models and not to train the models or adjust the hyperparameters.

3.2 NLU Models

We conduct the empirical study on five different kinds of pre-trained model: BERT, RoBERTa,

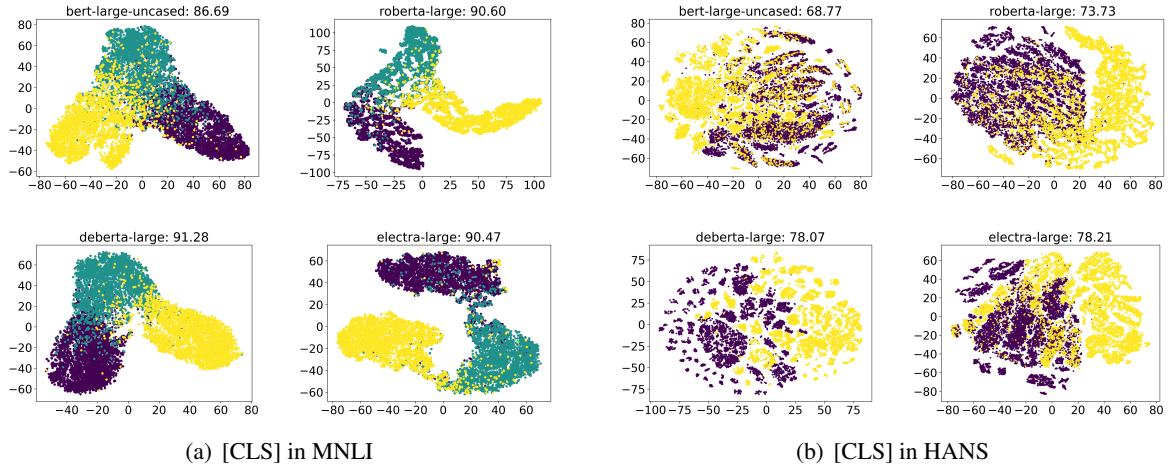


Figure 1: The t-SNE visualization result of [CLS] embeddings from different pre-trained language models on MNLI and HANS. The words above figures are the selected model and corresponding accuracy.

BART, ELECTRA, and DeBERTa. These models are used as encoders for the NLU models, which can provide contextual word embeddings. We use the corresponding pre-trained models and fine-tuned models from Huggingface Transformers². The input fed into the encoders of the NLU models is a pair of concatenated premise sentences and hypothesis sentences, which are separated by a special [SEP] token. Then, we obtain sentence pair representations through these models, which are the [CLS] embeddings from encoders. These representations are fed into the classification head (that is, the classifier) of the NLU models. Here, the classification head takes a simple architecture, two linear layers with the activation function.

3.3 Implementation Details

For all pre-trained models, we fine-tune the models without debiasing methods and with the example re-weighting debiasing method for 3 epochs. We find that the models converge slowly when fine-tuning the models with the product-of-experts debiasing method. Thus, we follow He et al. (2019) to fine-tune longer, i.e., 6 epochs. We use AdamW optimizer (Loshchilov and Hutter, 2019) with the default learning rate 5×10^{-5} , where the betas are set as [0.9, 0.999] and the L2 weight decay is set to 0.01. We set the batch size to 32 and warmup ratio to 0.1. All experiments are run with 3 random seeds and the average values are reported, which are completed on the work station with 2 Nvidia 2080Ti GPUs.

²<https://huggingface.co/models>

4 Experimental Analysis

In this section, we first use the visualization technique to investigate the separability of representations (§4.1). Then we investigate the linear separability of representations using DIRECTPROBE (§4.2). Furthermore, we propose an explanation for the difficulty in improving robustness by computing the similarity of representations (§4.3). Finally, we analyze the effect of representations and classifiers on robustness by considering debiasing methods as a probing tool (§4.4).

4.1 Visualization of Representations

To investigate the semantic representation space learned by the model, we extract embeddings of the special classification token [CLS] in the final hidden state and visualize them using t-SNE (Van der Maaten and Hinton, 2008). Figures 1(a) and 1(b) show the visualization results for MNLI and HANS, respectively. We show only the results of BERT_{large}, RoBERTa_{large}, DeBERTa_{large} and ELECTRA_{large}. Appendix B includes the results for more models. Based on the results, we find that the better performance the model achieves, the clearer the boundaries the representation forms. We believe that this is the reason why high performance is achieved with only a simple two-layer MLP network as the classifier.

In addition, we visualize the mean embeddings of all tokens in the final hidden state, which are abbreviated as [MEAN]. Note that the models are not trained with [MEAN]. The results of MNLI and HANS are shown in Figures 2(a) and 2(b), respec-

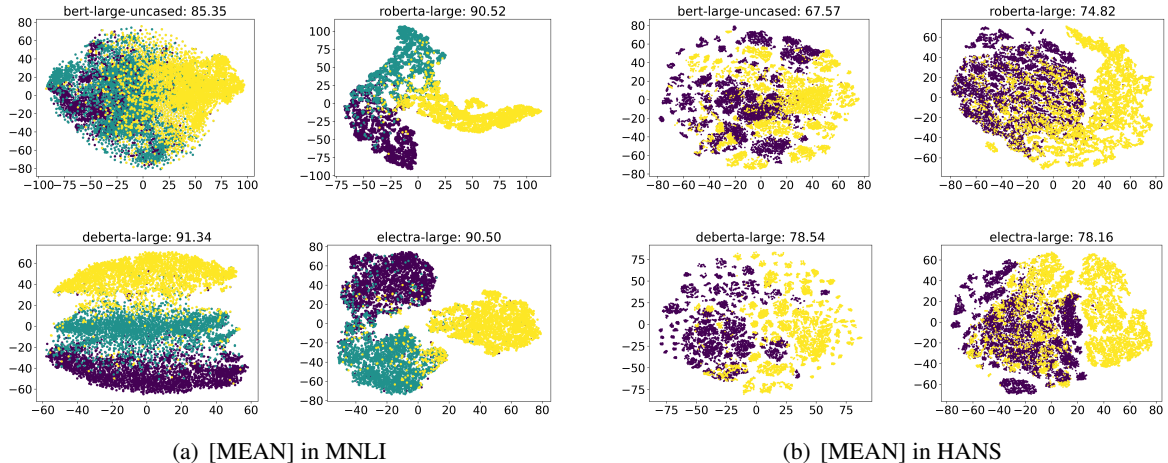


Figure 2: The t-SNE visualization result of [MEAN] embeddings from different pre-trained language models on MNL and HANS. The words above figures are the selected model and corresponding accuracy with [MEAN].

338 tively. Similarly, we show the results for more models
 339 in Appendix B. We observe that the [MEAN]
 340 embeddings can also form clearer boundaries as
 341 the models achieve better performance, despite the
 342 [CLS] embeddings are fed into classifiers. Mean-
 343 while, we freeze the encoders of the models trained
 344 with [CLS] and feed [MEAN] into the classifiers,
 345 the performance of which is close to the original
 346 one. Finally, we show the similarity between [CLS]
 347 and [MEAN] in Appendix C.

348 4.2 Linear Separability of Representations

349 In Figure 1(b), we discover the gap of improve-
 350 ment between the visualization effect and the per-
 351 formance, which may be due to the limitations of
 352 visualization technology. This leads us to introduce
 353 another method to quantify the quality of represen-
 354 tations. Therefore, we apply a probing technique
 355 based on the idea of clustering—DIRECTPROBE to
 356 study the geometric structure of representations.

357 We select five common pre-trained models to
 358 examine the geometric structure of representations
 359 after fine-tuning. Table 1 shows the results that
 360 contain base and large versions corresponding to
 361 selected models. We discover that the better per-
 362 formance the model achieves, the fewer clusters
 363 the representation is divided into, i.e., the represen-
 364 tation has higher linear separability. In particular,
 365 there are some models whose representations are
 366 divided into two and three clusters on HANS and
 367 MNL, respectively (i.e., equaling the number of
 368 label categories), which suggests that all examples
 369 with the same label are in one cluster and there are
 370 no overlaps between each cluster. However, these

Models	MNL		HANS	
	#clusters	Acc	#clusters	Acc
BERT _{base}	27	84.25	4	65.01
RoBERTa _{base}	5	88.10	3	69.53
DeBERTa _{base}	4	88.75	2	76.61
BART _{distill}	3	89.56	2	67.37
ELECTRA _{base}	4	88.77	2	76.84
BERT _{large}	5	86.69	2	68.77
RoBERTa _{large}	3	90.60	2	73.73
DeBERTa _{large}	3	91.28	2	78.07
BART _{large}	3	90.16	2	72.88
ELECTRA _{large}	3	90.47	2	78.21

Table 1: The number of clusters and corresponding accuracy from selected pre-trained language models on MNL and HANS. There is an around 20% generalization gap between MNL and HANS.

models achieve about 90% accuracy on MNL but
 only no more than 80% accuracy on HANS, which
 is not consistent with the linear separability of the
 representation. Thus, we assume that **the low ro-
 bustness of the NLU models is not due to the
 inseparability of representations.**

377 4.3 Similarity of Representations

378 In §4.1 and §4.2, we analyze representations with
 379 the t-SNE and DIRECTPROBE technique, respec-
 380 tively. However, what puzzles us is why the rep-
 381 resentation is linearly separable, while the perfor-
 382 mance of PLMs is not perfect. To study that, we
 383 investigate the cosine similarity of representations
 384 to find the reason for the flawed performance. In
 385 practice, we investigate the similarity between cluster
 386 centers within MNL or HANS and between

Models	M-E H-E	M-E H-N	H-E H-N	Acc
BERT _{base}	0.9376	0.8128	0.8597	65.01
RoBERTa _{base}	0.9562	0.8714	0.8131	69.53
DeBERTa _{base}	0.9733	0.7196	0.6921	76.61
BART _{distill}	0.9138	0.8067	0.9064	67.37
ELECTRA _{base}	0.9656	0.6550	0.6911	76.84
BERT _{large}	0.9573	0.7859	0.8511	68.77
RoBERTa _{large}	0.9286	0.7451	0.6810	73.73
DeBERTa _{large}	0.9545	0.6721	0.6912	78.07
BART _{large}	0.9145	0.7126	0.8467	72.88
ELECTRA _{large}	0.9560	0.5650	0.5891	78.21

Table 2: The similarity of [CLS] between two centers of selected embeddings and the accuracy on HANS. M-E indicates MNLi-Entailment; H-E indicates HANS-Entailment; H-N indicates HANS-Not-Entailment. It suggests that the representations with opposite semantics are similar in the semantic representation space.

MNLi and HANS. Each cluster includes all [CLS] embeddings with the same label on one dataset. The cluster center is defined as the average of all [CLS] embeddings in the cluster.

First, we compute the cosine similarity between MNLi-Entailment and HANS-Entailment/HANS-Not-Entailment as shown in Table 2. Ideally, the cosine similarity between MNLi-Entailment and HANS-Entailment is close to +1, and the cosine similarity between MNLi-Entailment and HANS-Not-Entailment is close to -1. However, the latter is not supported by our experiments, as shown in Table 2. It is even larger than +0.5, suggesting that the representations of HANS-Entailment and HANS-Not-Entailment are similar to those of MNLi-Entailment in the semantic representation space. Then, we compute the cosine similarity between HANS-Entailment and HANS-Not-Entailment as Table 2 shows, which should be close to -1. In fact, it is large than +0.5, suggesting that the representations of HANS-Entailment and HANS-Not-Entailment are similar, despite the representations from most of the selected models are linearly separable as Table 1 shows. These results mean that the encoders fail to distinguish Not-Entailment examples where the heuristics fail from Entailment examples well. In Appendix D, the gap of performance between Entailment and Not-Entailment on HANS confirm that.

Motivated by the above finding, we compute the Spearman correlation coefficient between similarity and accuracy and find that the accuracy is significantly inversely associated with the similarity between MNLi/HANS-Entailment and HANS-Not-Entailment. The correlation coefficients are

-0.8788 and -0.8909 with a p-value less than 0.05. We suppose that **the significant similarity in semantic representation between MNLi/HANS-Entailment and HANS-Not-Entailment is the main reason why it is difficult to improve the performance on HANS**. Finally, we show the other similarity between cluster centers with different labels in Appendix E.

4.4 Analysis of Debiasing Methods

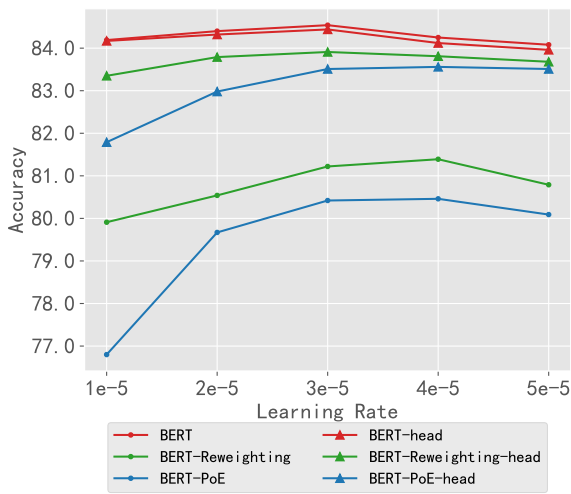
Typically, the poor performance of models on OOD examples is attributed to the fact that models only capture shortcut features (i.e., spurious features) but not robustness features (i.e. task-relevant features). Consequently, many debiasing methods are proposed to make models pay more attention to robustness features to improve the performance on OOD examples. In this work, we make detailed analyses of the impact of debiasing methods on fine-tuning models. For pre-trained language models, we select BERT and RoBERTa.

Based on the fact that the NLI task is considered as a classification task, the models for the NLI task can be divided into encoders and classifiers. Generally speaking, the encoder is from PLMs, and the classifier is the MLP network. Based on this architecture of the NLI models, we intend to study how debiasing methods work on encoders and classifiers, respectively. In practice, we apply a two-phase training strategy. Specifically, we first fine-tune models consisting of encoders and classifiers, then fix encoders and only retrain classifiers.

Figures 3 and 4 show the results on MNLi and HANS for BERT and RoBERTa, respectively. We take the learning rate from $1 * 10^{-5}$ to $5 * 10^{-5}$ to investigate the effect of the learning rate on the convergence of models. We discover that by fixing the encoder of fine-tuned models without debiasing methods (i.e. raw fine-tuned models) and only retraining the classifier with debiasing methods, the performance of models is improved on HANS for BERT and RoBERTa. Especially for BERT, the retraining strategy with debiasing methods achieves better performance than fine-tuning the whole model using debiasing methods both on MNLi and on HANS. Thus, we assume that **debiasing methods distort the representation to some degree** for BERT, which is similar to the finding of Mendelson and Belinkov (2021). Meanwhile, this strategy mitigates the degradation of performance on MNLi and further improves the performance on HANS. Noting that retraining the classifiers of

Models	Learning Rate				
	1e-5	2e-5	3e-5	4e-5	5e-5
BERT-ReW	79.91	80.54	81.22	81.39	80.79
BERT-ReW-head-self	79.92	80.61	81.29	81.39	80.90
BERT-PoE	76.80	79.67	80.42	80.46	80.09
BERT-PoE-head-self	77.06	79.90	80.58	80.61	80.32
RoBERTa-ReW	85.04	85.49	85.33	84.90	84.52
RoBERTa-ReW-head-self	85.18	85.41	85.29	84.77	84.44
RoBERTa-PoE	84.20	84.63	84.80	84.66	84.07
RoBERTa-PoE-head-self	84.21	84.78	84.97	84.60	84.14

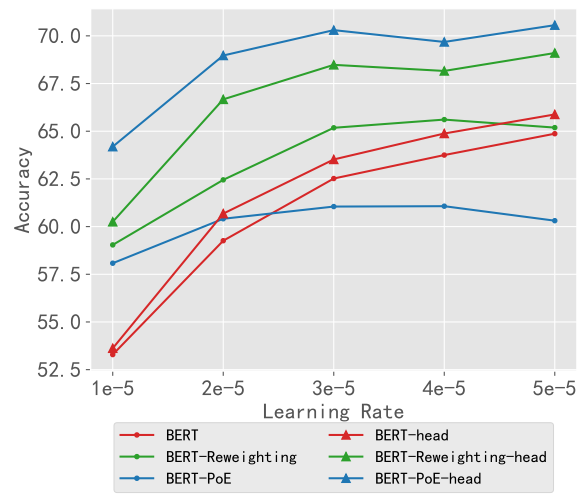
Table 3: The results of the model trained with the debiasing method and retraining the classifier of that with the same debiasing method for MNLI. There is no obvious change on performance before and after retraining.



(a) BERT-MNLI

Models	Learning Rate				
	1e-5	2e-5	3e-5	4e-5	5e-5
BERT-ReW	59.04	62.45	65.18	65.61	65.19
BERT-ReW-head-self	59.35	62.70	65.17	65.70	65.31
BERT-PoE	58.08	60.41	61.05	61.07	60.31
BERT-PoE-head-self	58.01	60.12	60.95	61.02	60.10
RoBERTa-ReW	75.66	77.69	77.74	76.40	75.90
RoBERTa-ReW-head-self	75.75	77.72	77.70	76.39	75.82
RoBERTa-PoE	77.29	78.50	78.29	77.64	76.27
RoBERTa-PoE-head-self	77.04	78.37	78.17	77.72	76.19

Table 4: The results of the model trained with the debiasing method and retraining the classifier of that with the same debiasing method for HANS. There is no obvious change on performance before and after retraining.



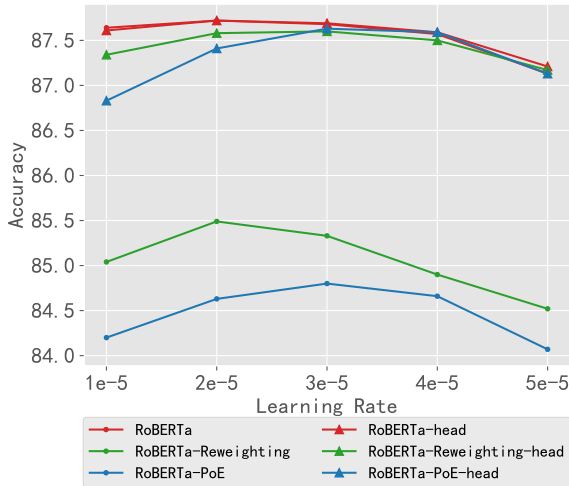
(b) BERT-HANS

Figure 3: (a) and (b) indicate the results of BERT on MNLI and HANS, respectively. The lines with dots: fine-tune models with or without debiasing methods. The lines with triangles: retrain classifiers using encoders from fine-tuned models with or without debiasing methods.

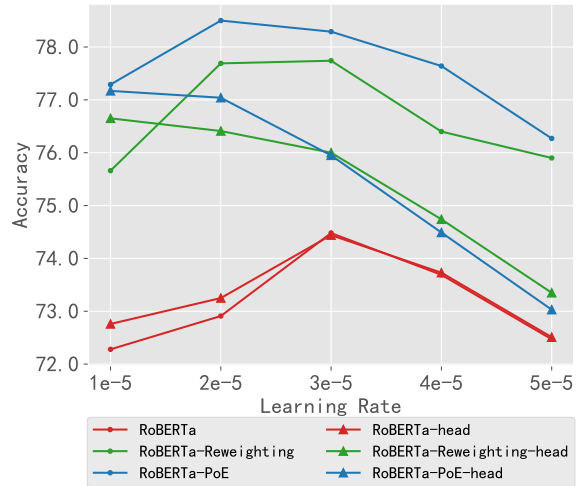
raw fine-tuned models without debiasing methods does not achieve an obvious effect on the performance, we find **classifiers play a significant role in the shortcut learning behavior of PLMs**. We design the comparative experiments to further clarify whether the performance improvement derives from the decoupled training strategy. Tables 3 and 4 show the results for MNLI and HANS, respectively. We discover that the decoupled training strategy also does not achieve an obvious effect for fine-tuned models with debiasing methods. This suggests that the performance improvement is derived from not the decoupled training strategy but the classifiers optimized by debiasing methods, and that the models fine-tuned with debiasing methods are not limited in performance by the classifiers.

Based on the above results, we suppose that a

decoupled retraining strategy with debiasing methods can be considered as a probing tool, which is used to measure the shortcut learning behavior of PLMs from representations or classifiers. The performance gap between the raw fine-tuned model and the model with the classifier retrained using the debiasing methods measures the extent of shortcut learning behavior from classifiers. The gap of the performance between the model with the classifier retrained using the debiasing methods and the whole fine-tuned model with the debiasing methods measures the extent of shortcut learning behavior from representations. Through this probing tool, we can better understand how representations and classifiers affect the robustness of models.



(a) RoBERTa-MNLI



(b) RoBERTa-HANS

Figure 4: (a) and (b) indicate the results of RoBERTa on MNLI and HANS, respectively. The lines with dots: fine-tune models with or without debiasing methods. The lines with triangles: retrain classifiers using encoders from fine-tuned models with or without debiasing methods.

5 Related Work

Recently, the shortcut learning behavior for the language task is revealed in previous work (Niven and Kao, 2019; Mudrakarta et al., 2018; Geirhos et al., 2020). For the NLI task, the shortcut learning behavior in models is often investigated using challenge datasets (Jia and Liang, 2017; Naik et al., 2018; Glockner et al., 2018; McCoy et al., 2019). To mitigate this behavior, we can use advanced pre-trained language models to obtain better representations (Liu et al., 2019; Lewis et al., 2020; Clark et al., 2020b; He et al., 2021), or apply debiasing methods to fine-tune language models (Schuster et al., 2019; Clark et al., 2019; Utama et al., 2020b; Utama et al., 2020a).

There are lots of works that analyze and understand learned representations with probing techniques. For instance, Tenney et al. (2019), Hewitt et al. (2021) and Whitney et al. (2021) consider classifiers as probes. Meanwhile, Mimno and Thompson (2017), Ethayarajh (2019) and Zhou and Srikumar (2021) inspect the representations from a geometric perspective. There are also efforts to understand pre-trained representations (Chen et al., 2021; Li et al., 2021) and fine-tuned ones (Zhou and Srikumar, 2022) respectively. In contrast, we focus our analysis on biased features and classifiers, and study the role that the quality of representations and the capability of classifiers play in the robustness of models respectively.

6 Conclusion

In this work, we conduct an empirical study on how the robustness of language models is affected by encoders and classifiers, respectively.

On the one hand, we show that *the low robustness of language models is not primarily due to representations not being easily separable*. i) We find that several excellent models provide linearly separable representations, which suggests that classifiers limit the performance of models. ii) We find that the significantly high similarity between representations with opposite semantics from in-distribution and out-of-distribution datasets is a reason for the low robustness.

On the other hand, we show *the relative role of representations and classifiers in the low robustness of language models*. i) We find that debiasing methods do not always improve the quality of representations but rather improve the performance of models with optimal classifiers. ii) We find that the robustness of models depends not only on the low quality of representations, but also on the capability of classifiers, and their ratios vary for different architectures and fine-tuning processes.

Finally, we hope that the insights obtained from the empirical analysis will be beneficial to the community, allowing them to pay more attention to the important roles of classifiers for models and design better solutions to alleviate shortcut learning and improve the robustness of PLMs in NLU tasks.

565 Limitations

566 Despite our findings that both representations and
567 classifiers affect the robustness of models, we are
568 not successful in making use of that to further im-
569 prove the understanding of models for the language.
570 As a result, we plan to further research advanced
571 methods that are capable of optimizing encoders
572 and classifiers, respectively. Furthermore, the de-
573 signed experiments in our analysis focus only on
574 the NLI task in NLU tasks. Given the similarity
575 between the NLU tasks, it may be possible to ex-
576 trapolate the corresponding findings to other NLU
577 tasks. In the future, we will consider the following
578 NLU tasks and datasets: IMDB (Maas et al., 2011)
579 / Yelp (Zhang et al., 2015) for the sentiment clas-
580 sification task; QQP (Iyer et al., 2017) / TwitterP-
581 PDB(TPPDB) (Lan et al., 2017) for the paraphrase
582 identification task; FEVER (Thorne et al., 2018) /
583 FeverSymmetric (Schuster et al., 2019) for the fact
584 verification task.

585 Ethics Statement

586 This paper does not raise ethical concerns. This
587 study does not involve any human subjects, prac-
588 tices to data set releases, potentially harmful in-
589 sights, discrimination/bias/fairness concerns and
590 privacy and security issues.

591 References

- 592 Yonatan Belinkov. 2021. [Probing classifiers: Promises,](#)
593 [shortcomings, and alternatives.](#) *arXiv e-prints*, pages
594 arXiv–2102.
- 595 Boli Chen, Yao Fu, Guangwei Xu, Pengjun Xie,
596 Chuanqi Tan, Mosha Chen, and Liping Jing. 2021.
597 [Probing {bert} in hyperbolic spaces.](#) In *International*
598 *Conference on Learning Representations*.
- 599 Christopher Clark, Mark Yatskar, and Luke Zettlemoyer.
600 2019. [Don’t take the easy way out: Ensemble based](#)
601 [methods for avoiding known dataset biases.](#) In *Pro-*
602 *ceedings of the 2019 Conference on Empirical Meth-*
603 *ods in Natural Language Processing and the 9th In-*
604 *ternational Joint Conference on Natural Language*
605 *Processing (EMNLP-IJCNLP)*, pages 4069–4082,
606 Hong Kong, China. Association for Computational
607 Linguistics.
- 608 Christopher Clark, Mark Yatskar, and Luke Zettlemoyer.
609 2020a. [Learning to model and ignore dataset bias](#)
610 [with mixed capacity ensembles.](#) In *Findings of the*
611 *Association for Computational Linguistics: EMNLP*
612 *2020*, pages 3031–3045, Online. Association for
613 Computational Linguistics.

- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and
614 Christopher D. Manning. 2020b. [Electra: Pre-](#)
615 [training text encoders as discriminators rather than](#)
616 [generators.](#) In *International Conference on Learning*
617 *Representations*. 618
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
619 Kristina Toutanova. 2019. [BERT: Pre-training of](#)
620 [deep bidirectional transformers for language under-](#)
621 [standing.](#) In *Proceedings of the 2019 Conference of*
622 *the North American Chapter of the Association for*
623 *Computational Linguistics: Human Language Tech-*
624 *nologies, Volume 1 (Long and Short Papers)*, pages
625 4171–4186, Minneapolis, Minnesota. Association for
626 Computational Linguistics. 627
- Mengnan Du, Varun Manjunatha, Rajiv Jain, Ruchi
628 Deshpande, Franck Dernoncourt, Jiuxiang Gu, Tong
629 Sun, and Xia Hu. 2021. [Towards interpreting and](#)
630 [mitigating shortcut learning behavior of NLU models.](#)
631 In *Proceedings of the 2021 Conference of the North*
632 *American Chapter of the Association for Computa-*
633 *tional Linguistics: Human Language Technologies,*
634 *pages 915–929, Online. Association for Computa-*
635 *tional Linguistics.* 636
- Kawin Ethayarajh. 2019. [How contextual are contextu-](#)
637 [alized word representations? Comparing the geom-](#)
638 [etry of BERT, ELMo, and GPT-2 embeddings.](#) In
639 *Proceedings of the 2019 Conference on Empirical*
640 *Methods in Natural Language Processing and the*
641 *9th International Joint Conference on Natural Lan-*
642 *guage Processing (EMNLP-IJCNLP)*, pages 55–65,
643 Hong Kong, China. Association for Computational
644 Linguistics. 645
- Robert Geirhos, Jörn-Henrik Jacobsen, Claudio
646 Michaelis, Richard Zemel, Wieland Brendel,
647 Matthias Bethge, and Felix A Wichmann. 2020.
648 [Shortcut learning in deep neural networks.](#) *Nature*
649 *Machine Intelligence*, 2(11):665–673. 650
- Max Glockner, Vered Shwartz, and Yoav Goldberg.
651 2018. [Breaking NLI systems with sentences that](#)
652 [require simple lexical inferences.](#) In *Proceedings*
653 *of the 56th Annual Meeting of the Association for*
654 *Computational Linguistics (Volume 2: Short Papers)*,
655 pages 650–655, Melbourne, Australia. Association
656 for Computational Linguistics. 657
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy,
658 Roy Schwartz, Samuel Bowman, and Noah A. Smith.
659 2018. [Annotation artifacts in natural language infer-](#)
660 [ence data.](#) In *Proceedings of the 2018 Conference of*
661 *the North American Chapter of the Association for*
662 *Computational Linguistics: Human Language Tech-*
663 *nologies, Volume 2 (Short Papers)*, pages 107–112,
664 New Orleans, Louisiana. Association for Computa-
665 tional Linguistics. 666
- He He, Sheng Zha, and Haohan Wang. 2019. [Unlearn](#)
667 [dataset bias in natural language inference by fitting](#)
668 [the residual.](#) In *Proceedings of the 2nd Workshop on*
669 *Deep Learning Approaches for Low-Resource NLP*
670 *(DeepLo 2019)*, pages 132–142, Hong Kong, China.
671 Association for Computational Linguistics. 672

673	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. {DEBERTA}: {DECODING}-{enhanced} {bert} {with} {disentangled} {attention} . In <i>International Conference on Learning Representations</i> .	(Volume 1: Long Papers), pages 4215–4228, Online. Association for Computational Linguistics.	730 731
674			
675			
676			
677			
678	Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus) . <i>arXiv preprint arXiv:1606.08415</i> .		732
679			733
680			734
681	John Hewitt, Kawin Ethayarajh, Percy Liang, and Christopher Manning. 2021. Conditional probing: measuring usable information beyond a baseline . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 1626–1639, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.		735
682			736
683			
684			
685			
686			
687			
688	Shankar Iyer, Nikhil Dandekar, and Kornél Csernai. 2017. Quora question pairs. <i>First Quora Dataset Release: Question Pairs</i> .		737
689			738
690			739
691	Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 2021–2031, Copenhagen, Denmark. Association for Computational Linguistics.		740
692			741
693			742
694			743
695			744
696			745
697			746
698			747
699			
700	Rabeeh Karimi Mahabadi, Yonatan Belinkov, and James Henderson. 2020. End-to-end bias mitigation by modelling biases in corpora . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8706–8716, Online. Association for Computational Linguistics.		748
701			749
702			750
703			751
704			752
705			753
706			754
707			755
708			756
709			757
710			758
711			759
712			760
713			
714			
715			
716			
717			
718			
719			
720			
721			
722			
723			
724			
725			
726			
727			
728			
729			
730			
731			
732			
733			
734			
735			
736			
737			
738			
739			
740			
741			
742			
743			
744			
745			
746			
747			
748			
749			
750			
751			
752			
753			
754			
755			
756			
757			
758			
759			
760			
761			
762			
763			
764			
765			
766			
767			
768			
769			
770			
771			
772			
773			
774			
775			
776			
777			
778			
779			
780			
781			
782			
783			
784			
785			
786			
787			
788			
789			
790			
791			
792			
793			
794			
795			
796			
797			
798			
799			
800			

787	Christian S Perone, Roberto Silveira, and Thomas S	from unknown biases. In <i>Proceedings of the 2020</i>	845
788	Paula. 2018. Evaluation of sentence embeddings	<i>Conference on Empirical Methods in Natural Lan-</i>	846
789	in downstream and linguistic probing tasks . <i>arXiv</i>	<i>guage Processing (EMNLP)</i> , pages 7597–7610, On-	847
790	<i>preprint arXiv:1806.06259</i> .	line. Association for Computational Linguistics.	848
791	Yada Pruksachatkun, Jason Phang, Haokun Liu,	Laurens Van der Maaten and Geoffrey Hinton. 2008.	849
792	Phu Mon Htut, Xiaoyi Zhang, Richard Yuanzhe Pang,	Visualizing data using t-sne. <i>Journal of machine</i>	850
793	Clara Vania, Katharina Kann, and Samuel R. Bow-	<i>learning research</i> , 9(11).	851
794	man. 2020. Intermediate-task transfer learning with	William F Whitney, Min Jae Song, David Brandfon-	852
795	pretrained language models: When and why does it	brener, Jaan Altonaar, and Kyunghyun Cho. 2021.	853
796	work? In <i>Proceedings of the 58th Annual Meeting of</i>	Evaluating representations by the complexity of learn-	854
797	<i>the Association for Computational Linguistics</i> , pages	ing low-loss predictors . In <i>Neural Compression:</i>	855
798	5231–5247, Online. Association for Computational	<i>From Information Theory to Applications – Work-</i>	856
799	Linguistics.	<i>shop @ ICLR 2021</i> .	857
800	Ivan Sanchez, Jeff Mitchell, and Sebastian Riedel. 2018.	Adina Williams, Nikita Nangia, and Samuel Bowman.	858
801	Behavior analysis of NLI models: Uncovering the in-	2018. A broad-coverage challenge corpus for sen-	859
802	fluence of three factors on robustness . In <i>Proceedings</i>	tence understanding through inference . In <i>Proceed-</i>	860
803	<i>of the 2018 Conference of the North American Chap-</i>	<i>ings of the 2018 Conference of the North American</i>	861
804	<i>ter of the Association for Computational Linguistics:</i>	<i>Chapter of the Association for Computational Lin-</i>	862
805	<i>Human Language Technologies, Volume 1 (Long Pa-</i>	<i>guistics: Human Language Technologies, Volume</i>	863
806	<i>pers)</i> , pages 1975–1985, New Orleans, Louisiana.	<i>1 (Long Papers)</i> , pages 1112–1122, New Orleans,	864
807	Association for Computational Linguistics.	Louisiana. Association for Computational Linguis-	865
808	Victor Sanh, Thomas Wolf, Yonatan Belinkov, and	tics.	866
809	Alexander M Rush. 2021. Learning from others’	Ruibin Xiong, Yimeng Chen, Liang Pang, Xueqi Cheng,	867
810	mistakes: Avoiding dataset biases without model-	Zhi-Ming Ma, and Yanyan Lan. 2021. Uncertainty	868
811	ing them . In <i>International Conference on Learning</i>	calibration for ensemble-based debiasing methods .	869
812	<i>Representations</i> .	<i>Advances in Neural Information Processing Systems</i> ,	870
813	Tal Schuster, Darsh Shah, Yun Jie Serene Yeo, Daniel	34.	871
814	Roberto Filizzola Ortiz, Enrico Santus, and Regina	Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.	872
815	Barzilay. 2019. Towards debiasing fact verification	Character-level convolutional networks for text clas-	873
816	models . In <i>Proceedings of the 2019 Conference on</i>	sification . In <i>Advances in Neural Information Pro-</i>	874
817	<i>Empirical Methods in Natural Language Processing</i>	<i>cessing Systems</i> , volume 28. Curran Associates, Inc.	875
818	<i>and the 9th International Joint Conference on Natu-</i>	Yichu Zhou and Vivek Srikumar. 2021. DirectProbe:	876
819	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages	Studying representations without classifiers . In <i>Pro-</i>	877
820	3419–3425, Hong Kong, China. Association for Com-	<i>ceedings of the 2021 Conference of the North Amer-</i>	878
821	putational Linguistics.	<i>ican Chapter of the Association for Computational</i>	879
822	Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang,	<i>Linguistics: Human Language Technologies</i> , pages	880
823	Adam Poliak, R Thomas McCoy, Najoung Kim, Ben-	5070–5083, Online. Association for Computational	881
824	jamin Van Durme, Sam Bowman, Dipanjan Das, and	Linguistics.	882
825	Ellie Pavlick. 2019. What do you learn from con-	Yichu Zhou and Vivek Srikumar. 2022. A closer look	883
826	text? probing for sentence structure in contextualized	at how fine-tuning changes BERT . In <i>Proceedings</i>	884
827	word representations . In <i>International Conference</i>	<i>of the 60th Annual Meeting of the Association for</i>	885
828	<i>on Learning Representations</i> .	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	886
829	James Thorne, Andreas Vlachos, Oana Cocarascu,	pages 1046–1061, Dublin, Ireland. Association for	887
830	Christos Christodoulopoulos, and Arpit Mittal. 2018.	Computational Linguistics.	888
831	The fact extraction and VERification (FEVER)		
832	shared task . In <i>Proceedings of the First Workshop on</i>		
833	<i>Fact Extraction and VERification (FEVER)</i> , pages 1–		
834	9, Brussels, Belgium. Association for Computational		
835	Linguistics.		
836	Prasetya Ajie Utama, Nafise Sadat Moosavi, and Iryna		
837	Gurevych. 2020a. Mind the trade-off: Debiasing		
838	NLU models without degrading the in-distribution		
839	performance . In <i>Proceedings of the 58th Annual</i>		
840	<i>Meeting of the Association for Computational Lin-</i>		
841	<i>guistics</i> , pages 8717–8729, Online. Association for		
842	Computational Linguistics.		
843	Prasetya Ajie Utama, Nafise Sadat Moosavi, and Iryna		
844	Gurevych. 2020b. Towards debiasing NLU models		

A Implementation Details

Only-bias Model The only-bias model is trained on random 2,000 of examples for 3 epochs. We use AdamW optimizer with the default learning rate $5 * 10^{-5}$, where the betas are set as [0.9, 0.999] and the L2 weight decay is set to 0.01. The batch size is set to 32, and the warmup ratio is set to 0.1.

Main Model The hidden dimension of the classifiers is the same as the output of encoders (i.e., base version: 768; large version: 1024). The activation functions of the classifiers are the same as the setup of the Huggingface Transformers (i.e., Tanh or GELU (Hendrycks and Gimpel, 2016)): i) BERT, BART, and RoBERTa are Tanh; ii) DeBERTa and ELECTRA are GELU. The parameters of PLMs are shown as Table 5.

Retrain Classifiers with or without Debiasing Methods The parameters of the classifier are initialized by a normal distribution with the mean of 0.0 and the variance of 0.02. We use AdamW optimizer with the default learning rate $5 * 10^{-5}$, where the betas are set to [0.9, 0.999] and the L2 weight decay is set to 0.01. The batch size is set to 32 and the warmup ratio is set to 0.1. We retrain the classifiers for 3 epochs.

Models	Parameters
BERT _{base}	110M
BERT _{large}	340M
RoBERTa _{base}	125M
RoBERTa _{large}	355M
DeBERTa _{base}	134M
DeBERTa _{large}	390M
BART _{distill}	356M
BART _{large}	406M
ELECTRA _{base}	110M
ELECTRA _{large}	335M

Table 5: The parameters of pre-trained language models.

B Other Results of Visualization

Figures 5 and 6 show the other visualization results for [CLS] and [MEAN], respectively.

C Similarity between [CLS] and [MEAN]

To explore how [CLS] and [MEAN] are related in terms of robustness, we compute the cosine simi-

ilarity between [CLS] and [MEAN] on MNLI and HANS, respectively. Table 6 summarizes the results. We compare the change in similarity from base models to large ones and discover that the changes in similarity have different trends. When the trend of the change in similarity increases, we suppose that the model is likely to learn similar information. On the contrary, the model is likely to learn different information. Based on this observation, we conjecture that it is possible to improve the robustness of models by figuring out how the amount of information learned affects performance and introducing the information from [MEAN] as supervised signals while fine-tuning. Verifying or rejecting this conjecture requires further study.

Models	MNLI		HANS	
	Similarity	Acc	Similarity	Acc
BERT _{base}	0.7683	84.25	0.7441	65.01
BERT _{large}	0.6364	86.69	0.6146	68.77
RoBERTa _{base}	0.8224	88.10	0.7663	69.53
RoBERTa _{large}	0.9845	90.60	0.9914	73.73
DeBERTa _{base}	0.5718	88.75	0.5690	76.61
DeBERTa _{large}	0.3362	91.28	0.2371	78.07
BART _{distill}	0.6375	89.56	0.6903	67.37
BART _{large}	0.5989	90.16	0.6433	72.88
ELECTRA _{base}	0.5188	88.77	0.6048	76.84
ELECTRA _{large}	0.8461	90.47	0.9317	78.21

Table 6: The cosine similarity between [CLS] and [MEAN] and corresponding accuracy from selected pre-trained language models on MNLI and HANS.

Models	MNLI		HANS	
	[CLS]	[MEAN]	[CLS]	[MEAN]
BERT _{base}	84.25	84.17	65.01	65.20
RoBERTa _{base}	88.10	87.94	69.53	70.89
DeBERTa _{base}	88.75	88.89	76.61	78.12
BART _{distill}	89.56	88.86	67.37	63.20
ELECTRA _{base}	88.77	88.47	76.84	76.21
BERT _{large}	86.69	85.35	68.77	67.57
RoBERTa _{large}	90.60	90.52	73.73	74.82
DeBERTa _{large}	91.28	91.34	78.07	78.54
BART _{large}	90.16	88.95	71.88	71.34
ELECTRA _{large}	90.47	90.50	78.21	78.16

Table 7: The results of [CLS] and [MEAN] for MNLI and HANS.

D Results of HANS in detail

Table 10 shows the results of HANS in detail.

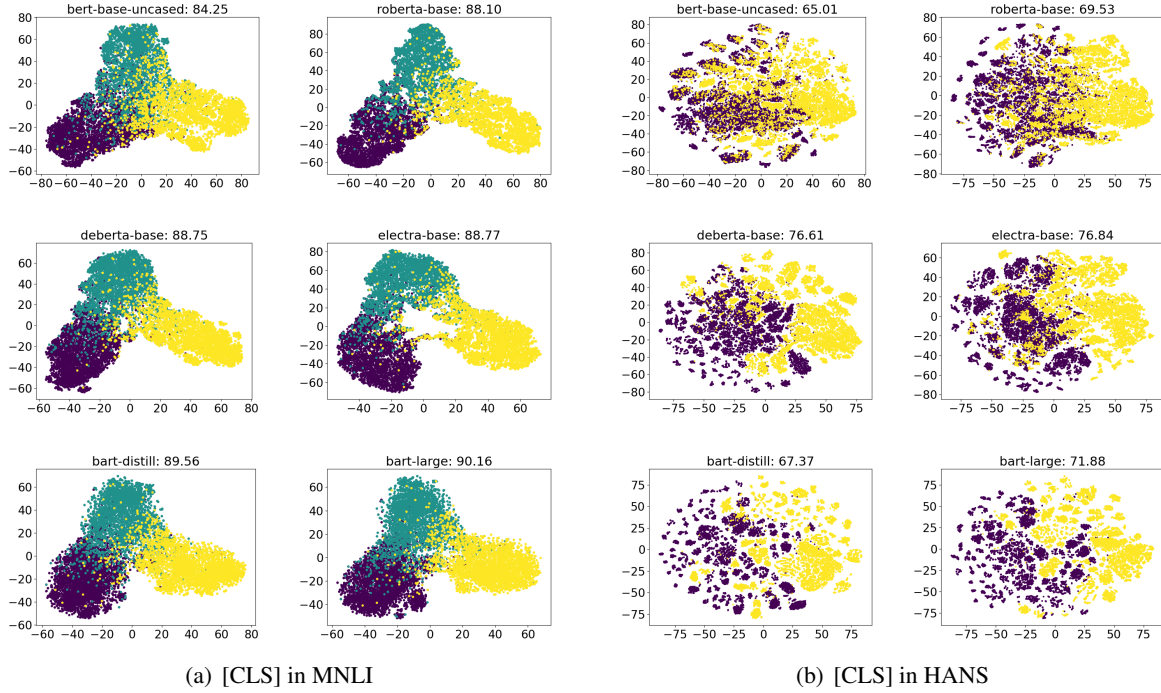


Figure 5: The t-SNE visualization result of [CLS] embeddings from different pre-trained language models on MNLI and HANS. The words above figures are the selected model and corresponding accuracy.

Models	Learning Rate				
	1e-5	2e-5	3e-5	4e-5	5e-5
BERT	84.19	84.40	84.54	84.25	84.08
BERT-ReW	79.91	80.54	81.22	81.39	80.79
BERT-PoE	76.80	79.67	80.42	80.46	80.09
BERT-head	84.17	84.32	84.44	84.12	83.96
BERT-ReW-head	83.35	83.79	83.91	83.81	83.68
BERT-PoE-head	81.79	82.98	83.51	83.56	83.51
BERT-ReW-head-self	79.92	80.61	81.29	81.39	80.90
BERT-PoE-head-self	77.06	79.90	80.58	80.61	80.32
RoBERTa	87.64	87.72	87.68	87.57	87.13
RoBERTa-ReW	85.04	85.49	85.33	84.90	84.52
RoBERTa-PoE	84.20	84.63	84.80	84.66	84.07
RoBERTa-head	87.61	87.72	87.69	87.59	87.21
RoBERTa-ReW-head	87.34	87.58	87.60	87.50	87.17
RoBERTa-PoE-head	86.83	87.41	87.63	87.59	87.13
RoBERTa-ReW-head-self	85.18	85.41	85.29	84.77	84.44
RoBERTa-PoE-head-self	84.21	84.78	84.97	84.60	84.14

Table 8: The complete results for MNLI.

E Similarity between Cluster Centers

Table 11 shows the complete similarity between cluster centers within MNLI or HANS and between MNLI and HANS.

F Complete Results

Tables 8 and 9 show the complete results of BERT and RoBERTa for MNLI and HANS, respectively.

Models	Learning Rate				
	1e-5	2e-5	3e-5	4e-5	5e-5
BERT	53.29	59.26	62.52	63.75	64.87
BERT-ReW	59.04	62.45	65.18	65.61	65.19
BERT-PoE	58.08	60.41	61.05	61.07	60.31
BERT-head	53.62	60.68	63.52	64.88	65.88
BERT-ReW-head	60.25	66.67	68.48	68.16	69.10
BERT-PoE-head	64.19	68.97	70.30	69.68	70.56
BERT-ReW-head-self	59.35	62.70	65.17	65.70	65.31
BERT-PoE-head-self	58.01	60.12	60.95	61.02	60.10
RoBERTa	72.28	72.91	74.48	73.69	72.47
RoBERTa-ReW	75.66	77.69	77.74	76.40	75.90
RoBERTa-PoE	77.29	78.50	78.29	77.64	76.27
RoBERTa-head	72.76	73.25	74.44	73.73	72.51
RoBERTa-ReW-head	76.65	76.41	76.00	74.74	73.35
RoBERTa-PoE-head	77.17	77.04	75.95	74.49	73.03
RoBERTa-ReW-head-self	75.75	77.72	77.70	76.39	75.82
RoBERTa-PoE-head-self	77.04	78.37	78.17	77.72	76.19

Table 9: The complete results for HANS.

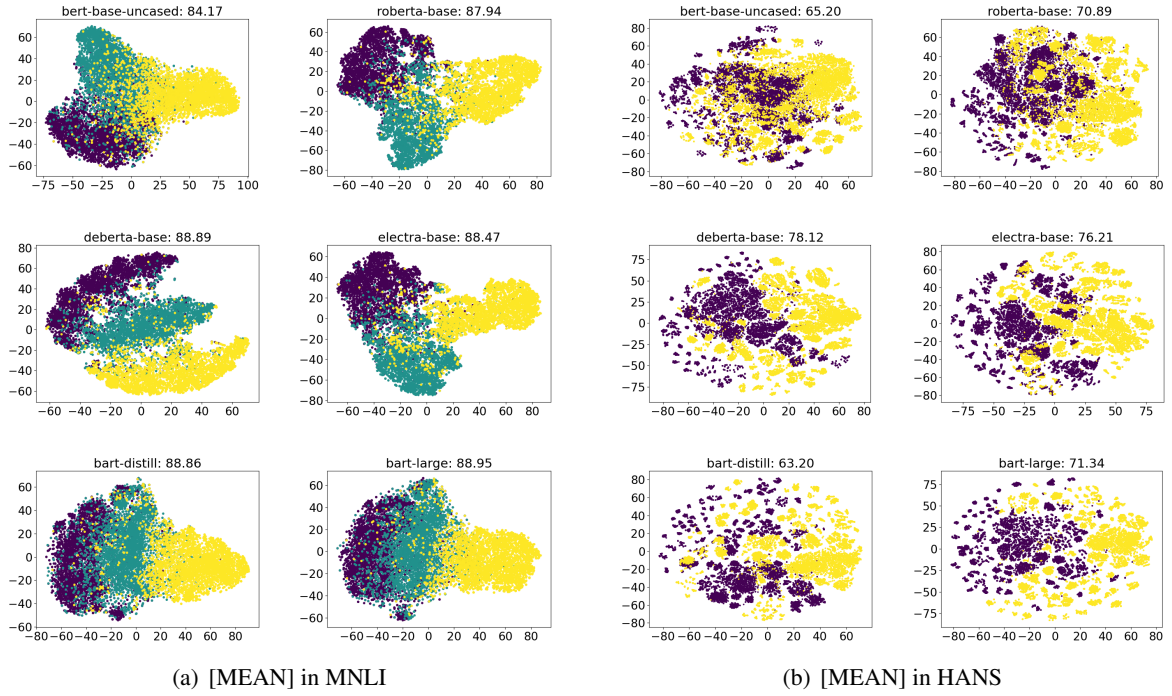


Figure 6: The t-SNE visualization result of [MEAN] embeddings from different pre-trained language models on MNL and HANS. The words above figures are the selected model and corresponding accuracy.

Models	HANS	HANS-Entailment	HANS-Not-Entailment	Entailment Category			Non-Entailment Category		
				Overlap	Subsequence	Constituent	Overlap	Subsequence	Constituent
BERT _{base}	65.01	98.97	31.05	97.54	99.64	99.74	59.10	12.12	21.94
RoBERTa _{base}	69.53	99.39	39.66	98.96	99.98	99.24	66.02	19.72	33.24
DeBERTa _{base}	76.61	99.21	54.01	97.82	100.0	99.80	95.60	30.70	35.72
BART _{distill}	67.37	99.25	35.49	98.32	99.72	99.70	69.18	16.36	20.94
ELECTRA _{base}	76.84	99.59	54.09	98.90	99.94	99.94	95.92	27.98	38.38
BERT _{large}	68.77	94.85	42.69	88.22	97.60	98.74	74.90	22.62	30.56
RoBERTa _{large}	73.73	99.63	47.83	99.98	100.00	98.92	90.52	34.82	18.14
DeBERTa _{large}	78.07	99.86	56.27	99.74	100.00	99.84	95.00	33.28	40.54
BART _{large}	71.88	99.53	44.24	99.02	99.76	99.80	80.76	27.32	24.64
ELECTRA _{large}	78.21	99.84	56.58	99.52	100.00	100.00	93.04	37.24	39.46

Table 10: The results of HANS in detail.

Models	M-EI-H-E	M-EI-H-N	M-NI-H-E	M-NI-H-N	M-CI-H-E	M-CI-H-N	M-EIM-N	M-EIM-C	M-NIM-C	H-EI-H-N	Acc
BERT _{base}	0.9376	0.8128	-0.0725	0.0470	-0.1906	0.2225	0.1744	-0.2330	0.1160	0.8597	65.01
RoBERTa _{base}	0.9562	0.8714	0.1196	0.3833	-0.1000	0.3951	0.3640	-0.0354	0.2031	0.8131	69.53
DeBERTa _{base}	0.9733	0.7196	0.1833	0.4205	-0.1263	0.4709	0.2963	-0.1289	0.1063	0.6921	76.61
BART _{distill}	0.9138	0.8067	0.1695	0.2850	0.1791	0.4585	0.3210	0.1178	0.2796	0.9064	67.37
ELECTRA _{base}	0.9656	0.6550	-0.1790	-0.0977	-0.3068	0.3648	0.0150	-0.3798	-0.1271	0.6911	76.84
BERT _{large}	0.9573	0.7859	0.0926	0.3132	0.1037	0.4651	0.2327	0.0019	0.2867	0.8511	68.77
RoBERTa _{large}	0.9286	0.7451	-0.1081	0.0972	-0.0796	0.4204	0.1048	-0.1673	-0.1164	0.6810	73.73
DeBERTa _{large}	0.9545	0.6721	0.1204	0.3466	-0.0546	0.5147	0.2880	-0.1025	0.0611	0.6912	78.07
BART _{large}	0.9145	0.7126	0.0698	0.1178	0.1050	0.4824	0.1828	-0.0084	0.0782	0.8467	72.88
ELECTRA _{large}	0.9560	0.5650	-0.3017	-0.1654	-0.1751	0.5167	-0.1574	-0.2974	-0.2505	0.5891	78.21

Table 11: The similarity of [CLS] between two centers of selected embeddings and the accuracy on HANS. M-E indicates MNL-Entailment; M-N indicates MNL-Neutral; M-C indicates MNL-Contradiction; H-E indicates HANS-Entailment; H-N indicates HANS-Not-Entailment.