

# Prompt-Based Bias Calibration for Better Zero/Few-Shot Learning of Language Models

Anonymous ACL submission

## Abstract

Prompt learning is susceptible to intrinsic bias present in pre-trained language models (LMs), resulting in sub-optimal performance of prompt-based zero/few-shot learning. In this work, we propose a *null-input prompting* method to calibrate intrinsic bias encoded in pre-trained LMs. Different from prior efforts that address intrinsic bias primarily for social fairness and often involve excessive computational cost, our objective is to explore enhancing LMs' performance in downstream zero/few-shot learning while emphasizing the efficiency of intrinsic bias calibration. Specifically, we leverage a diverse set of auto-selected null-meaning inputs generated from GPT-4 to probe intrinsic bias of pre-trained LMs. Utilizing the bias-reflected probability distribution, we formulate a distribution disparity loss for bias calibration, where we exclusively update bias parameters (0.1% of total parameters) of LMs towards equal probability distribution. Experimental results show that the calibration promotes an equitable starting point for LMs while preserving language modeling abilities. Across a wide range of datasets, including sentiment analysis and topic classification, our method significantly improves zero/few-shot learning performance of LMs for both in-context learning and prompt-based fine-tuning (on average 9% and 2%, respectively).

## 1 Introduction

The advent of GPT models (Radford et al., 2019; Brown et al., 2020) has catalyzed the transformative prompt-learning paradigm. The innovative approach of "pre-train, prompt, and predict" (Schick and Schütze, 2021a; Liu et al., 2023) facilitates fast adaptation of pre-trained language models (LMs) in learning various tasks and empowering LMs' strong zero/few-shot learning abilities (Schick and Schütze, 2021b; Gao et al., 2021).

Due to the susceptibility to bias ingrained in pre-trained LMs, prompt learning tends to make

biased predictions toward some specific answers, thereby impacting the performance of prompt-based zero/few-shot learning (Zhao et al., 2021; Han et al., 2023). To mitigate this issue and improve LM performance, Zhao et al. (2021) and Holtzman et al. (2022) propose to reweigh LM output probabilities. Han et al. (2023) explores calibrating decision boundaries. While these research has demonstrated substantial improvements, they are primarily designed for in-context learning with frozen pre-trained LMs, leading to two main limitations: (1) They may be not effective in task-specific fine-tuning scenario (Jian et al., 2022). Note, however, prompt-based fine-tuning has shown performance improvements over in-context learning (Gao et al., 2021; Logan IV et al., 2022). It is particularly important for relatively small-sized LMs. (2) The intrinsic bias encoded in pre-trained LMs persists since these research focuses on *output calibration* and does not modify LMs.

To address these limitations, we investigate the potential for enhancing the performance of LMs as zero/few-shot learners in classification tasks by *calibrating intrinsic bias* of pre-trained LMs. This exploration extends to various prompt-learning scenarios: in-context learning and prompt-based fine-tuning. Prior approaches to mitigate intrinsic bias primarily focus on achieving social fairness, and often require laborious corpora augmentation and costly re-training (Huang et al., 2020; Kaneko and Bollegala, 2021; Solaiman and Dennison, 2021; Li et al., 2023a). To improve efficiency in both data generation and model updates, we propose leveraging auto-generated *null-meaning inputs* to prompt pre-trained LMs for intrinsic bias probing, and subsequently updating only *bias parameters*  $B_{LM}$  of LMs for bias calibration. Null-meaning inputs are essentially normal text devoid of meaningful content or sentiment. Unlike numerical-zero inputs, they maintain the contextual framework of prompts, ensuring the proper functioning of contex-

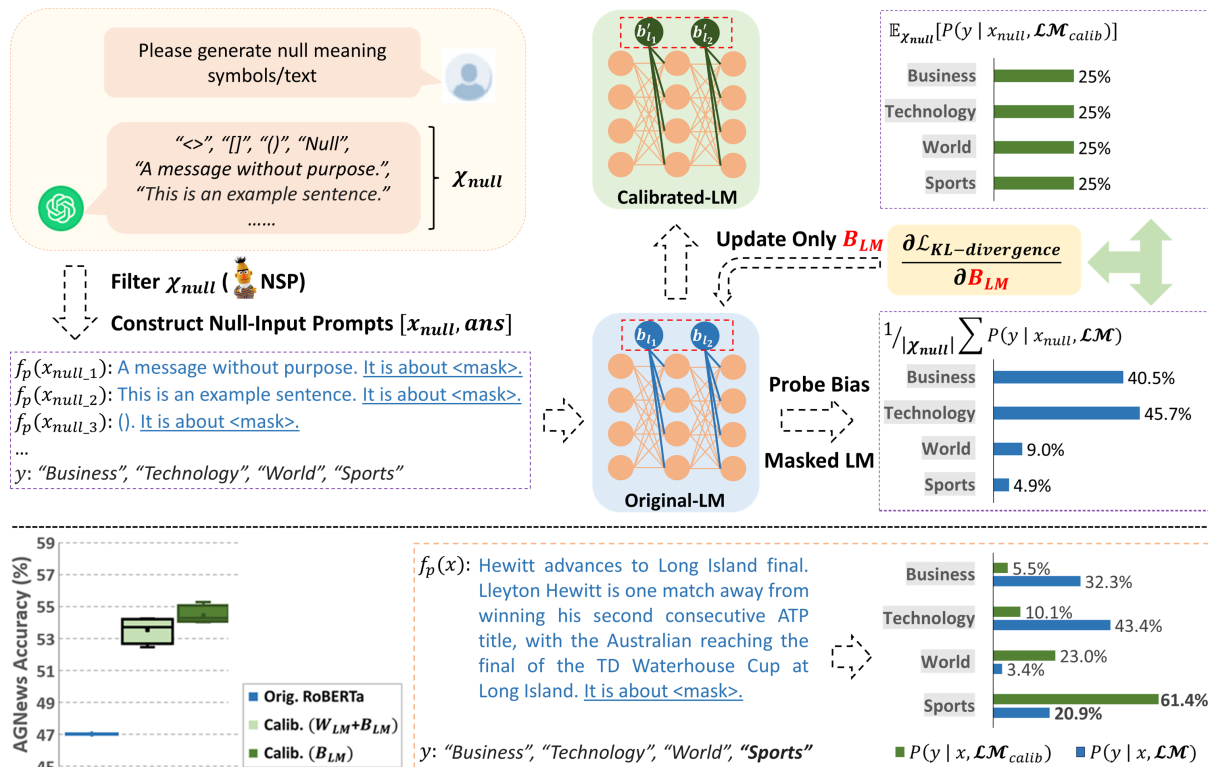


Figure 1: We demonstrate our calibration method significantly improves classification performance of pre-trained LM. **Upper**: The pipeline of proposed null-input prompting method for intrinsic bias calibration targeting AGNews task (Zhang et al., 2015). **Lower left**: Performance comparison of zero-shot in-context learning using: original LM (Orig. RoBERTa); calibrated (Calib.) LM with full model updates ( $W_{LM} + B_{LM}$ ); calibrated LM with only  $B_{LM}$  updates. **Lower right**: Case study illustrating that LM makes correct prediction after intrinsic bias calibration.

085 tual LMs. Our motivation stems from the expectation that bias-calibrated models should produce uniform probabilities across all categories if the input in a prompt delivers null information (Zhao et al., 2021).  $B_{LM}$  functions as offsets in neural networks, and strategically updating only  $B_{LM}$  could potentially counteract intrinsic bias of pre-trained models, achieving higher efficiency (updating  $\sim 0.1\%$  parameters of entire LM). The approach promotes an equitable starting point, and we expect that the light model updates preserve pre-trained models’ language modeling abilities while maintaining the focus on bias calibration, ultimately making LMs better zero/few-shot learners.

099 The pipeline of our calibration method is illustrated in Figure 1. We use Masked LMs (RoBERTa Liu et al., 2019) for zero/few-shot learning since they generally produce competitive performance in classification tasks and their moderate size facilitates combining prompting with fine-tuning (Gao et al., 2021; Liu et al., 2023). First, we utilize GPT-4 API to automatically generate diverse null-meaning inputs  $\mathcal{X}_{null}$  including symbols, words, phrases, and sentences. This generation process is

109 downstream task-agnostic. By concatenating each null-meaning input  $x_{null}$  with an answer format aligned with the downstream task, we construct null-input prompts (similar to Zhao et al., 2021), e.g., “An empty sentence. It is about <mask>.”. For better cohesive integration of the “null” information into the prompts, we additionally devise a filtering strategy to select  $x_{null}$ , to which the answer format *ans* exhibits relatively strong Next Sentence Prediction (NSP) correlation (Devlin et al., 2019). Next, we update  $B_{LM}$  with null-input prompts to calibrate intrinsic bias. Given the absence of task-relevant information in these prompts, the anticipated outcome in the parameter updating process is a convergence towards equal output probabilities for each label word. We formulate a customized Kullback–Leibler (KL) divergence loss for gradient descent on  $B_{LM}$  to minimize the distribution disparity. Finally, bias-calibrated LMs are applied in downstream prompt-based zero/few-shot learning following Gao et al. (2021).

The main contributions of our work are:

- We introduce a null-input prompting method for calibrating intrinsic bias of pre-trained

Masked LMs, aiming for better prompt-based zero/few-shot classification performance.

- Our method integrates two key aspects for efficient bias calibration: auto-construction of null-input prompts and updating only bias parameters of LMs. The calibration promotes a fair starting point for LMs while preserving language modeling abilities.
- Extensive experiments on eight classification datasets with four prompt-learning approaches show that our method significantly improves LMs’ zero/few-shot performance, and outperforms output-calibration methods.

## 2 Related Work

**Impact of intrinsic bias on downstream LM performance.** Intrinsic bias in pre-trained LMs stems from imbalances present in extensive pre-training corpora. Higher frequency of specific terms in those corpora could lead to *common token bias* (Zhao et al., 2021). Additionally, frequent co-occurrence of certain terms with specific sentiment in pre-training could introduce *association bias* (Cao et al., 2022). Because of those intrinsic bias, prompt-based predictions by pre-trained LMs are prone to bias towards some specific answers, resulting in sub-optimal performance in downstream tasks (Zhao et al., 2021; Han et al., 2023).

**Mitigating strategies.** Research has focused on counteracting the bias solely at the output prediction stage, without modifying pre-trained LMs. For example, Zhao et al. (2021) introduces contextual calibration and Holtzman et al. (2022) presents Domain Conditional Pointwise Mutual Information to reweigh answer scores. Min et al. (2022) explores computing the probability of the input conditioned on the label. Han et al. (2023) proposes to calibrate decision boundaries. However, these studies mainly demonstrate their effectiveness for in-context learning using frozen pre-trained LMs, without addressing the intrinsic bias encoded in the LMs. Other research on mitigating intrinsic bias primarily targets removing social bias (Dinan et al., 2020; Huang et al., 2020; Cheng et al., 2021; Zhou et al., 2023), often employing costly data augmentation and re-training, and as a by-product, degrades language modeling abilities (Meade et al., 2022).

Efficiently calibrating intrinsic bias in pre-trained LMs for enhancing downstream zero/few-shot learning performance is an open research problem. We introduce a parameter-efficient intrinsic-

bias calibration method leveraging automatically constructed null-input prompts, which significantly improves zero/few-shot learning of LMs.

**Parameter-efficient fine-tuning (PEFT) for downstream tasks.** It has been demonstrated that fine-tuning a very small portion of model parameters can achieve performance on par with fine-tuning the entire set of parameters. People propose integrating small, trainable adapter modules between model layers (Bapna and Firat, 2019; Houlsby et al., 2019), coupled with further optimization using low-rank adaptations (LoRA) (Hu et al., 2021). Some other research focuses on prompt tuning (Lester et al., 2021; Li and Liang, 2021; Gu et al., 2022; Guo et al., 2022) which only tunes continuous prompt embeddings for efficiently adapting pre-trained LMs to downstream tasks.

Our method provides a unique perspective of enhancing LM performance on downstream tasks through efficient intrinsic-bias calibration. We update only bias parameters of pre-trained LMs with null-input prompts in calibration. Contrary to adapters and LoRA which would need sufficient labeled data to learn new matrices, we do not introduce new matrices to pre-trained LMs, preserving LMs’ few-shot learning capabilities. Moreover, our approach does not necessarily require target-domain data (whether labeled or unlabeled), enabling fully unsupervised deployment, particularly advantageous for zero-shot setting.

## 3 Null-Input Prompting for Intrinsic Bias Calibration

### 3.1 Task Formulation

Let  $\mathcal{LM}$  be a pre-trained Masked LM. Verbalizer  $V(\cdot)$  maps label  $y$  to vocabulary token. Prompt function  $f_p(\cdot)$  modifies original input  $x_{in}$  into cloze-style prompt containing one  $\langle \text{mask} \rangle$  token to be predicted. The output representation  $\mathbf{h}_{\langle \text{mask} \rangle}$  of the  $\langle \text{mask} \rangle$  token is acquired from the last encoder layer after forwarding the prompt to the LM. Following Gao et al. (2021), the probability prediction of each class  $y \in \mathcal{Y}$  is formulated as:

$$\begin{aligned} P(y | x_{in}, \mathcal{LM}) &= P(V(y) | f_p(x_{in}), \mathcal{LM}) \\ &= \frac{\exp(\text{index}_{V(y)}(\mathbf{W}_{lm\_head} \cdot \mathbf{h}_{\langle \text{mask} \rangle}))}{\sum_{j=1}^{|\mathcal{Y}|} \exp(\text{index}_{V(y_j)}(\mathbf{W}_{lm\_head} \cdot \mathbf{h}_{\langle \text{mask} \rangle}))}, \end{aligned} \quad (1)$$

where  $\mathbf{W}_{lm\_head}$  is the pre-trained *masked language modeling head* weight matrix, and  $\text{index}_{V(y)}$  se-

lects the logits corresponding to the label words based on their index in LM token list.

One can probe intrinsic bias encoded in pre-trained LM by replacing  $x_{in}$  with null-meaning input  $x_{null} \in \mathcal{X}_{null}$  (Zhao et al., 2021).  $\mathcal{X}_{null}$  represents a set of  $x_{null}$  and we will elaborate their generation and selection in § 4. As shown by the blue bars in the upper part of Figure 1, while null-meaning inputs essentially provide no task-relevant prior information, the mean output probability associated with different labels  $\bar{P}_{\mathcal{X}_{null}}(y | x_{null}, \mathcal{LM})$  may exhibit significant difference attributed to model’s intrinsic bias. Ideally, for bias-calibrated LM  $\mathcal{LM}_{calib}$ , the expectation of output distribution conditioned on null-meaning inputs should be uniform across all label words, i.e.,

$$\mathbb{E}_{\mathcal{X}_{null}} [P(y | x_{null}, \mathcal{LM}_{calib}; \forall y \in \mathcal{Y})] = \frac{1}{|\mathcal{Y}|}. \quad (2)$$

We aim to calibrate intrinsic bias by updating LM to minimize this distribution disparity which we quantify using differentiable KL divergence as:

$$\begin{aligned} & D_{KL} (U(\mathcal{Y}) || \bar{P}_{\mathcal{X}_{null}}(\mathcal{Y})) \\ &= \sum_{y \in \mathcal{Y}} \left( 1/|\mathcal{Y}| \cdot \log \frac{1/|\mathcal{Y}|}{\bar{P}_{\mathcal{X}_{null}}(y)} \right) \\ &= \log(1/|\mathcal{Y}|) - (1/|\mathcal{Y}|) \cdot \sum_{y \in \mathcal{Y}} \log \bar{P}_{\mathcal{X}_{null}}(y), \quad (3) \end{aligned}$$

where  $U(\mathcal{Y})$  denotes uniform probability distribution and  $\bar{P}_{\mathcal{X}_{null}}(y)$  represents the simplified form of  $\bar{P}_{\mathcal{X}_{null}}(y | x_{null}, \mathcal{LM})$ .

### 3.2 Update Only Bias Parameters

While intrinsic bias may be encoded across various parts of pre-trained LMs, one question arises: is it essential to update the entire model, or is there a more efficient alternative that can achieve comparable effectiveness in intrinsic bias calibration? We propose to only update bias parameters  $\mathbf{B}_{LM}$ , with the following rationale: (i)  $\mathbf{B}_{LM}$  constitutes less than 0.1% of total LM parameters, offering significant memory and computation cost saving compared to updating entire LM. (ii) Weight parameters  $\mathbf{W}_{LM}$ <sup>1</sup> may carry crucial pre-existing knowledge for language modeling, which risks impairment with a full model update (Meade et al., 2022).  $\mathbf{B}_{LM}$ , often overlooked in LM research, serves as offsets in DNN layers. Strategic updates may counteract

<sup>1</sup> $\mathbf{W}_{LM}$  also includes embedding parameters in our context.

intrinsic bias while potentially preserving language modeling abilities. (iii) Empirical research on efficient fine-tuning has demonstrated the important role of bias parameters in LMs (Ben Zaken et al., 2022; Logan IV et al., 2022).

We update  $\mathbf{B}_{LM}$  using gradient descent to minimize the dissimilarity between output probability distribution from the LM conditioned on null-meaning inputs and uniform probability distribution  $U(\mathcal{Y})$ . We formulate a customized KL divergence loss  $\mathcal{L}$ , including both divergence of individual null-input’s output distribution  $P_i(\mathcal{Y})$  with respect to  $U(\mathcal{Y})$ , and batch-averaged distribution  $\bar{P}_N(\mathcal{Y})$  with respect to  $U(\mathcal{Y})$ , as:

$$\begin{aligned} \mathcal{L} = & \frac{1}{N} \sum_{i=1}^N D_{KL} (U(\mathcal{Y}) || P_i(\mathcal{Y})) \\ & + D_{KL} (U(\mathcal{Y}) || \bar{P}_N(\mathcal{Y})), \quad (4) \end{aligned}$$

where  $N$  is the batch size of null-meaning inputs. Incorporating the second term in the loss function promotes calibration stability and aligns with the objective of Equation 2.

### 3.3 Early Stopping of Calibration

We aim to obtain LM with improved zero/few-shot performance at the calibration stopping point. An overly calibrated model may simply produce uniform probability predictions regardless of input information. To avoid this, we develop specialized early stopping strategies depending on whether the downstream task is zero-shot or few-shot.

**For zero-shot downstream tasks.** Determining the calibration stopping point for optimal zero-shot learning performance is challenging due to the absence of labeled data for validation during calibration. To discern the patterns of a good stopping point, we first conduct empirical experiments by validating LM zero-shot performance on the entire test dataset after each calibration batch (consisting of  $N$  null-meaning inputs) across different calibration learning rates (Figure 7 in Appendix A). As shown in Figure 2, with optimal calibration learning rate, model performance exhibits significant improvements in the first one/few calibration batches with low variance, and then starts to degrade and becomes unstable. The low performance and instability at the calibration tail confirm our assumption on the detrimental effects of excessive calibration on LM’s modeling abilities. Notably, calibration with only one batch of null inputs (indicated by the red vertical line in Figure 2) delivers



consistent and significant improvement compared to the original LM (although might not be the best improvement). Therefore, for enhancing LM zero-shot performance, we directly adopt the *One-batch Calibration* as the early stopping criterion.

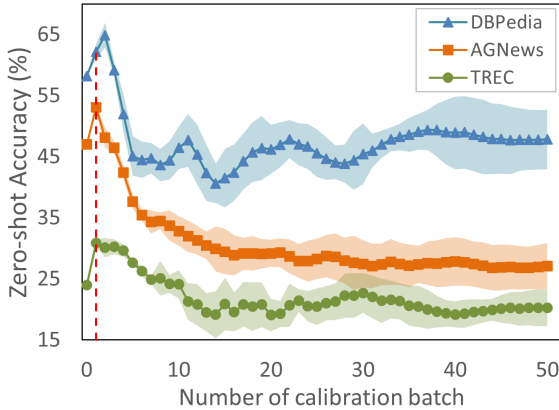


Figure 2: Empirical experiments show the impact of calibration on zero-shot learning performance as the number of calibration batches increases (batch size is 32). The intersections of the curves and red vertical line signify the outcomes of the first calibration batch.

**For few-shot downstream tasks.** With the acquisition of a few labeled downstream data, the previous challenge of lacking validation for determining the stopping point in the calibration process is alleviated. We utilize the small amount of labeled data as validation dataset  $\mathcal{D}_{\text{val}}^{\text{calib}}$  to set a stopping criterion for calibration. Additionally, we take into account above-mentioned empirical findings that, for some tasks, stopping at one batch of calibration yields optimal LM performance. Relying on the limited size of  $\mathcal{D}_{\text{val}}^{\text{calib}}$  might fail to identify such stopping points. To this effect, we store both  $LM_{\text{calib}}^{\text{one\_batch}}$  (obtained from one-batch stopping) and  $LM_{\text{calib}}^{\text{val}}$  (obtained from validation-based stopping) for downstream few-shot leaning tasks. Since  $LM_{\text{calib}}^{\text{one\_batch}}$  is stored in the process of obtaining  $LM_{\text{calib}}^{\text{val}}$ , this will not result in additional computation overhead. Memory overhead is minimal, as it only requires storing an additional set of updated bias parameters.

We summarize our method for intrinsic bias calibration in Algorithm 1 (Appendix A).

## 4 Auto-Construct Null-Input Prompt

### 4.1 Generate Null-Meaning Input

We employ null-meaning inputs to probe the intrinsic bias of pre-trained LMs, and then use those bias-reflected outputs to calibrate the LMs. Crafting a diverse set of null-meaning inputs  $\mathcal{X}_{\text{null}}$  for

an averaged output helps prevent overfitting to sub-optimal instances, thereby contributing to the effectiveness of calibration. To enable cost-effective acquisition of various null-meaning data, we utilize GPT-4 API for automatic generation with instructions such as "Please generate null meaning symbols, words, phrases, and sentences, in total <Number>.". This process is task-agnostic, generating data that contains null information with respect to any downstream task. Note that null information is not equivalent to neutral sentiment, as it carries no inherent meaning or contextual sentiment implications. We further validate this through t-SNE (van der Maaten and Hinton, 2008) visualization in Appendix A Figure 6.

Generated null-meaning input $x_{\text{null}}$	$P_{\text{NSP}}(x_{\text{null}}, \text{ans})$
<i>This is an example sentence.</i>	0.9996
<i>A message without purpose.</i>	0.9979
<i>Words without message.</i>	0.9809
<i>123abc</i>	0.0267
@#\$\$%^&*()-_+=+[]{} ////////////////	0.0145
////////////////	0.0008

Table 1: Some examples of generated null-mean inputs. In this case, "It is about <mask>." is used as the answer format *ans*. The green/yellow numbers represent higher/lower NSP probabilities.

### 4.2 Select $x_{\text{null}}$ and Build Null-Input Prompt

We construct null-input prompt  $f_p(x_{\text{null}})$  by concatenating the generated null-meaning input with an answer format *ans*. For consistency, the answer format (e.g., "It is <mask>.") is the same as the one intended for use in the downstream task. Some examples are shown in the upper part of Figure 1.

To pursue better cohesive integration of the "null" information into the prompts, we prioritize the null-meaning inputs, with which the answer format exhibits higher Next Sentence Prediction (NSP) probability (Devlin et al., 2019). Specifically, after we generate a large set of null-meaning inputs  $\{x_{\text{null}_1}, x_{\text{null}_2}, \dots, x_{\text{null}_k}\}$  and the answer format *ans* is selected, we employ BERT-large model (Devlin et al., 2019) to predict NSP  $P_{\text{NSP}}(x_{\text{null}}, \text{ans})$  and sort null-meaning inputs by their probabilities. Table 1 shows some generated  $x_{\text{null}}$ , with which a specific answer format presents high/low NSP scores. After the sorting, we retain the top 80%  $x_{\text{null}}$  instances (800 in total), which maintains the diversity among the selected samples. We observed that null inputs

	In-context lrn no demo <sup>†</sup>			In-context lrn with demo			Prompt FT no demo			Prompt FT with demo		
	NoCal	OutCal	IntrCal	NoCal	OutCal	IntrCal	NoCal	OutCal	IntrCal	NoCal	OutCal	IntrCal
AGNews	47.0 <sub>0.0</sub>	54.3 <sub>1.0</sub>	<b>54.5</b> <sub>0.6</sub>	79.7 <sub>0.8</sub>	78.8 <sub>3.3</sub>	<b>82.4</b> <sub>0.9</sub>	<b>89.1</b> <sub>0.9</sub>	86.3 <sub>1.6</sub>	89.0 <sub>0.8</sub>	86.9 <sub>2.8</sub>	87.5 <sub>1.3</sub>	<b>89.3</b> <sub>0.9</sub>
DBPedia	58.2 <sub>0.0</sub>	54.1 <sub>1.9</sub>	<b>61.8</b> <sub>0.6</sub>	92.6 <sub>0.6</sub>	94.0 <sub>0.9</sub>	<b>94.8</b> <sub>0.7</sub>	98.2 <sub>1.3</sub>	99.0 <sub>0.5</sub>	<b>99.0</b> <sub>0.1</sub>	98.6 <sub>0.3</sub>	98.5 <sub>0.2</sub>	<b>98.9</b> <sub>0.3</sub>
TREC	24.0 <sub>0.0</sub>	29.4 <sub>2.1</sub>	<b>31.1</b> <sub>0.5</sub>	48.3 <sub>1.4</sub>	42.5 <sub>3.4</sub>	<b>48.6</b> <sub>2.2</sub>	85.0 <sub>7.4</sub>	82.2 <sub>2.0</sub>	<b>89.3</b> <sub>4.5</sub>	87.6 <sub>2.5</sub>	74.2 <sub>4.0</sub>	<b>89.7</b> <sub>1.0</sub>
Subj	50.8 <sub>0.0</sub>	<b>64.0</b> <sub>2.7</sub>	62.7 <sub>0.8</sub>	47.2 <sub>0.2</sub>	55.0 <sub>1.3</sub>	<b>63.5</b> <sub>2.3</sub>	91.2 <sub>0.9</sub>	88.2 <sub>2.5</sub>	<b>93.2</b> <sub>1.2</sub>	91.4 <sub>3.3</sub>	93.0 <sub>0.8</sub>	<b>94.3</b> <sub>0.2</sub>
SST-5	31.5 <sub>0.0</sub>	33.0 <sub>2.1</sub>	<b>37.5</b> <sub>0.4</sub>	34.4 <sub>1.7</sub>	31.2 <sub>2.6</sub>	<b>36.6</b> <sub>1.0</sub>	47.8 <sub>4.6</sub>	45.3 <sub>2.8</sub>	<b>49.9</b> <sub>2.7</sub>	47.1 <sub>1.9</sub>	42.6 <sub>4.0</sub>	<b>50.0</b> <sub>1.7</sub>
Laptop	54.6 <sub>0.0</sub>	58.3 <sub>2.5</sub>	<b>59.6</b> <sub>1.9</sub>	50.8 <sub>1.0</sub>	65.1 <sub>2.7</sub>	<b>67.4</b> <sub>1.7</sub>	74.3 <sub>1.4</sub>	74.3 <sub>1.6</sub>	<b>74.9</b> <sub>2.9</sub>	76.8 <sub>1.0</sub>	75.6 <sub>1.4</sub>	<b>78.7</b> <sub>1.4</sub>
Restaurant	68.6 <sub>0.0</sub>	72.0 <sub>4.9</sub>	<b>72.8</b> <sub>1.6</sub>	69.8 <sub>1.1</sub>	<b>74.3</b> <sub>1.6</sub>	74.0 <sub>1.0</sub>	79.7 <sub>2.2</sub>	79.0 <sub>1.0</sub>	<b>82.0</b> <sub>0.9</sub>	78.4 <sub>4.9</sub>	79.0 <sub>5.5</sub>	<b>79.8</b> <sub>4.5</sub>
Twitter	19.7 <sub>0.0</sub>	43.4 <sub>4.1</sub>	<b>51.7</b> <sub>0.4</sub>	21.0 <sub>0.5</sub>	40.7 <sub>5.4</sub>	<b>49.4</b> <sub>2.7</sub>	51.7 <sub>2.9</sub>	44.1 <sub>3.9</sub>	<b>57.0</b> <sub>4.2</sub>	57.7 <sub>2.8</sub>	50.3 <sub>4.2</sub>	<b>59.3</b> <sub>2.3</sub>
Average	44.3	51.1	<b>54.0</b>	55.5	60.2	<b>64.6</b>	77.1	74.8	<b>79.3</b>	78.1	75.1	<b>80.0</b>

Table 2: Result comparisons among NoCal (LM-BFF Gao et al., 2021; no calibration), OutCal (output calibration) and IntrCal (ours; intrinsic-bias calibrated LM) using RoBERTa-large. We report the mean and standard deviation of performance in 8 classification datasets with 4 prompt-learning methods. "In-context lrn" refers to in-context learning and "Prompt FT" refers to prompt-based fine-tuning. "with/no demo" denotes incorporating/not incorporating demonstrations in prompts. In-context lrn no demo<sup>†</sup> is zero-shot learning, while the other three are few-shot learning.

with lower NSP scores are typically randomly-combined alphabet letters and symbols. These samples may have minimal occurrences in pre-training corpora. The low NSP scores can be attributed to RoBERTa’s lack of comprehension of their meanings in context. Their representations extracted by LM might have high variance, which might impact the stability and effectiveness of calibration. We show calibration with  $x_{\text{null}}$  selection strategy further improves LM performance in § 5.2 Table 3.

## 5 Experiments

We conduct extensive experiments on 8 English datasets, including sentiment analysis and topic classification.<sup>2</sup> They consist of 5 sentence-level datasets potentially impacted by *common token bias*: AGNews (Zhang et al., 2015), DBPedia (Lehmann et al., 2015), TREC (Voorhees and Tice, 2000), Subj (Pang and Lee, 2004), SST-5 (Socher et al., 2013) and 3 aspect-level sentiment analysis datasets likely subject to *association bias*: Restaurant and Laptop reviews from SemEval 2014 Task (Pontiki et al., 2014), Twitter (Dong et al., 2014). For aspect-level datasets, the task is to predict sentiments associated with the marked aspects in each sentence. More details are in Appendix A Table 7.

### 5.1 Evaluation Protocol

We evaluate the effectiveness of our intrinsic-bias calibration method on enhancing Masked LMs

<sup>2</sup>We mainly focus on single-sentence tasks, which aligns with the use of single-sentence null inputs for calibration. The alignment may enhance calibration effectiveness. We also experiment on sentence-pair tasks in Appendix B Table 15 and demonstrate better performance after calibration.

zero/few-shot learning performance with 4 prompt learning methods: in-context learning and prompt-based fine-tuning, both with and without demonstration. We follow the prompt-based fine-tuning and demonstration method of Gao et al. (2021).

We conduct calibration with 5 different seeds, and for the few-shot setting, we randomly sample 5 different groups of training and validation sets ( $K$  samples per class). We report the mean and standard deviation of LM performance. For the 5 sentence-level classification tasks, we use *accuracy* as the metric. For the 3 aspect-level classification tasks, because of the imbalance in test set, we use *weighted  $F_1$*  for a balanced evaluation. Details of calibration and prompt learning are in Appendix A.

We present our main results using RoBERTa-large, and  $K = 16$  for few-shot setting. Results of using RoBERTa-base,  $K = \{2, 4, 8\}$ , and different prompt templates are in Appendix B (Table 10, Table 11 and Figure 8).

### 5.2 Main Results

In Table 2, we compare our results of **IntrCal** (intrinsic bias calibration) with reproduced results of: (1) **NoCal**: No calibration. Use LM-BFF (Gao et al., 2021) to compute  $P(y | x_{\text{in}})$  for predictions. (2) **OutCal**: Output calibration. OutCal computes  $\frac{P(y | x_{\text{in}})}{P(y | x_{\text{domain}})}$  instead of  $P(y | x_{\text{in}})$  to counteract surface form competition and bias (Zhao et al., 2021; Holtzman et al., 2022). Note that OutCal was originally demonstrated for in-context learning with GPT models, while here, we apply the method in Masked LMs for fair comparisons.

In addition to NoCal and OutCal, we compare

our results with those reproduced from *NoisyTune* (Wu et al., 2022), *NSP-BERT* (Sun et al., 2022) and *Perplexion* (Lu et al., 2023), as detailed in Appendix B.1 (Table 8, 9). The superior performance further validates the effectiveness of our method.

**In-context learning results.** OutCal has significantly improved LM zero/few-shot performance compared to NoCal. Our method (IntrCal) further outperforms OutCal by a large margin: 2.9% and 8.3% absolute in zero-shot learning & 4.4% and 8.7% absolute in few-shot learning, in terms of average and best-case improvement. This demonstrates the advantages of intrinsic bias calibration over attempting to counteract bias solely at the output. Moreover, OutCal exhibits higher variance in performance due to its sensitivity to human-crafted domain-relevant strings  $x_{\text{domain}}$ . Using certain  $x_{\text{domain}}$  instances may not accurately capture the bias of LMs, resulting in under-calibration or over-calibration and leading to the high variance. In our approach, we use a large set of auto-generated and selected  $x_{\text{null}}$  as the training set for bias calibration. This mitigates the impact of sub-optimal samples and enhances calibration robustness, contributing to more stable and reliable performance.

**Prompt-based fine-tuning results.** This method fine-tunes all LM parameters utilizing limited labeled data by minimizing the cross-entropy loss based on Equation 1. It greatly raises LM performance compared to in-context learning and sets up a strong baseline (i.e., NoCal). OutCal fails to surpass NoCal. We speculate that OutCal’s limitation lies in its exclusive focus on offsetting bias at the output and lack of interaction with the interior of LM. This appears to impede OutCal from adapting effectively to the intricate dynamics of LM after prompt-based fine-tuning, leading to some counterproductive calibrations. In contrast, IntrCal (ours) with the aim of intrinsic bias calibration achieves superior performance with absolute gains of maximum 5.3% and average 2% compared to NoCal.

The output representations of <mask> token for label word predictions are visualized by t-SNE in Figure 3. On the left, samples from the two categories are almost mixed together, indicating that the original LM tends to bias toward one class prediction. In contrast, the right visualization demonstrates improved separability after *One-batch Calibration* (§ 3.3), which explains the significant performance enhancement achieved by our intrinsic-bias calibration method.

	In-context lrn no demo		Prompt FT no demo	
	UnSel. $x_{\text{null}}$	Sel. $x_{\text{null}}$	UnSel. $x_{\text{null}}$	Sel. $x_{\text{null}}$
AGNews	53.1 <sub>0.6</sub>	<b>54.5</b> <sub>0.6</sub>	87.8 <sub>1.7</sub>	<b>89.0</b> <sub>0.8</sub>
DBPedia	<b>62.1</b> <sub>1.2</sub>	61.8 <sub>0.6</sub>	98.7 <sub>0.2</sub>	<b>99.0</b> <sub>0.1</sub>
TREC	30.9 <sub>0.6</sub>	<b>31.1</b> <sub>0.5</sub>	88.5 <sub>3.5</sub>	<b>89.3</b> <sub>4.5</sub>
Subj	60.5 <sub>3.2</sub>	<b>62.7</b> <sub>0.8</sub>	92.8 <sub>1.6</sub>	<b>93.2</b> <sub>1.2</sub>
SST-5	35.5 <sub>1.7</sub>	<b>37.5</b> <sub>0.4</sub>	48.7 <sub>4.2</sub>	<b>49.9</b> <sub>2.7</sub>

Table 3: Benefits from null-meaning input  $x_{\text{null}}$  selection strategy (§ 4.2). *UnSel.* signifies using all GPT-generated  $x_{\text{null}}$  in calibration, while *Sel.* denotes selecting top  $x_{\text{null}}$  based on the sorting of  $P_{\text{nsp}}(x_{\text{null}}, \text{ans})$ .

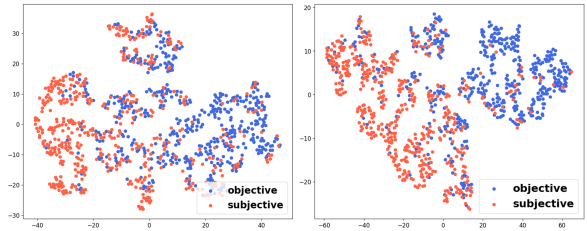


Figure 3: t-SNE visualization for output representations of <mask> token. **Left** is obtained from original LM; **Right** is obtained from the LM after *One-batch Calibration*. Two colors denote the two classes in Subj task.

	In-context lrn no demo		Prompt FT no demo	
	$W_{LM} + B_{LM}$	$B_{LM}$	$W_{LM} + B_{LM}$	$B_{LM}$
AGNews	53.5 <sub>0.8</sub>	<b>54.5</b> <sub>0.6</sub>	<b>89.3</b> <sub>0.8</sub>	89.0 <sub>0.8</sub>
DBPedia	<b>63.2</b> <sub>0.9</sub>	61.8 <sub>0.6</sub>	99.0 <sub>0.5</sub>	<b>99.0</b> <sub>0.1</sub>
TREC	<b>31.3</b> <sub>0.8</sub>	31.1 <sub>0.5</sub>	87.6 <sub>2.8</sub>	<b>89.3</b> <sub>4.5</sub>
Subj	53.3 <sub>0.6</sub>	<b>62.7</b> <sub>0.8</sub>	<b>93.7</b> <sub>0.6</sub>	93.2 <sub>1.2</sub>
SST-5	33.5 <sub>0.4</sub>	<b>37.5</b> <sub>0.4</sub>	49.4 <sub>0.7</sub>	<b>49.9</b> <sub>2.7</sub>
Laptop	58.2 <sub>0.8</sub>	<b>59.6</b> <sub>1.9</sub>	<b>78.1</b> <sub>1.3</sub>	74.9 <sub>2.9</sub>
Restaurant	70.7 <sub>1.8</sub>	<b>72.8</b> <sub>1.6</sub>	81.3 <sub>1.0</sub>	<b>82.0</b> <sub>0.9</sub>
Twitter	<b>51.8</b> <sub>0.7</sub>	51.7 <sub>0.4</sub>	55.7 <sub>2.3</sub>	<b>57.0</b> <sub>4.2</sub>
Average	51.9	<b>54.0</b>	79.3	79.3

Table 4: Performance comparisons between differently calibrated LMs.  $W_{LM} + B_{LM}$  updates entire LM in calibration while  $B_{LM}$  only updates bias parameters. Additional results of In-context lrn/Prompt FT *with demo* are in Appendix B Table 14.

### 5.3 Update Entire LM vs. Only Bias Parameters in Calibration

In Table 4, we evaluate the impact of updating entire LM ( $W_{LM} + B_{LM}$ ) during calibration on downstream task performance, as compared to only updating bias parameters ( $B_{LM}$ ). The optimal learning rate for updating entire LM is smaller (Appendix A Table 6). For in-context learning, the LM with only  $B_{LM}$  updates in calibration achieves better overall performance compared to the LM with entire

parameter updates, most likely attributed to better preserved language modeling abilities (Appendix B Table 12). For prompt-based fine-tuning, two differently calibrated LMs demonstrate comparable performance, as the impact of entire-parameter calibration on the modeling ability is mitigated through task-specific fine-tuning. Considering the significant saving in memory and computation, we recommend only updating  $B_{LM}$  in calibration.

## 5.4 Analysis

**How does intrinsic bias calibration impact downstream tasks?** Our method calibrates the intrinsic bias associated with a set of task-specific label words. In this section, we explore the impact of updating LM for task-specific bias calibration on other downstream task performance. Specifically, we take the LM calibrated for one task and evaluate its performance on the other tasks as shown in Figure 4. In general, intrinsic bias calibration for one task has a minimal adverse effect on other tasks’ performance (no more than 2% degradation) because of the light model updates, while remarkably enhancing LM performance on that specific task. Notably, there is consistent performance increase at bottom right, as these tasks are all sentiment classification sharing or including same label words.<sup>3</sup>

AGNews	47.0	+7.5	-2.0	+2.7	+1.6	+1.0	0.0	+0.5	+1.7
DBPedia	58.2	-1.4	+3.6	+3.0	+2.7	+2.0	+1.4	+2.4	+1.8
TREC	24.0	+1.4	-0.6	+7.1	+0.8	+1.4	+0.3	+0.9	-0.1
Subj	50.8	+0.6	-0.3	+0.2	+11.9	+0.1	-0.3	-0.3	-0.2
SST-5	31.5	-1.2	-0.6	+1.0	-0.7	+6.0	+4.5	+5.4	+5.2
Laptop	54.6	-1.5	-1.6	+1.4	-1.2	+4.3	+5.0	+3.7	+3.2
Restaurant	68.6	-0.3	-1.9	+1.7	-1.6	+4.4	+4.0	+4.2	+3.7
Twitter	19.7	-0.4	-0.8	+1.6	+0.7	+27.2	+29.0	+29.8	+32.0
	Baseline	AGNews	DBPedia	TREC	Subj	SST-5	Laptop	Restaurant	Twitter

Figure 4: Impact of calibration on downstream tasks shown through the changes with respect to baseline on each column. Each row shows the zero-shot performance of one task employing: *original LM* (first column; baseline), *task-specific calibrated LM* (diagonal), *other-task calibrated LM* (other places).

<sup>3</sup>For aspect-level datasets, better improvement is on the diagonals (task-specific calibration), indicating our method mitigates the impact of association bias (Appendix A).

**How does intrinsic bias calibration impact language modeling abilities?** We employ pseudo-perplexity (Salazar et al., 2020) to evaluate language modeling for Masked LM. Following each task-specific intrinsic bias calibration, we measure pseudo-perplexity and compare the results with original RoBERTa on WikiText-2, WikiText-103 (Merity et al., 2017), and LAMBADA dataset (Paperno et al., 2016). As shown in Table 5, language modeling abilities are largely preserved after calibration due to the minimal updates to the model.

	WT-2	WT-103	LAMBADA
Original RoBERTa	6.189	7.008	24.52
+ CALIBRATION			
for_AGNews	↑0.017 6.206	↑0.029 7.037	↑0.02 24.54
for_DBPedia	↑0.008 6.197	↑0.002 7.010	↓0.22 24.30
for_TREC	↓0.027 6.162	↓0.042 6.966	↓0.27 24.25
for_Subj	↓0.021 6.168	↓0.030 6.978	↑0.08 24.60
for_SST-5	↓0.031 6.158	↓0.039 6.969	↓0.18 24.34
for_Laptop	↑0.011 6.200	↑0.002 7.010	↓0.01 24.51
for_Restaurant	↑0.055 6.244	↑0.074 7.082	↑0.13 24.65
for_Twitter	↓0.029 6.160	↓0.037 6.971	↑0.05 24.57

Table 5: Pseudo-perplexities of *original RoBERTa* and *task-specific calibrated RoBERTa* on WikiText-2 (WT-2), WikiText-103 (WT-103) and LAMBADA. We use 2000 test samples of each dataset. An increase in values (highlighted in red) indicates a reduction in language modeling abilities after calibration.

## 6 Conclusion

In this work, we propose a null-input prompting method to calibrate the intrinsic bias of pre-trained Masked LMs, aiming to enhance zero/few-shot learning performance in classification tasks. Our method incorporates two key features for efficiency: (1) auto-construction of null-input prompts for bias probing, leveraging a diverse set of selected null-meaning inputs easily crafted from generative Large LM; (2) updating only bias parameters for bias calibration. Experimental results show that bias-calibrated LMs demonstrate significant performance improvement for both in-context learning and prompt-based fine-tuning, with average gains of 9% and 2%, respectively. Moreover, our method outperforms output-calibration approaches, highlighting the advantage of intrinsic bias calibration. We believe this work presents a new perspective of making LMs better zero/few-shot learners via intrinsic bias calibration. Additionally, the demonstrated significance of bias parameters could provide insights for future bias-related research.



## 7 Limitations

While our method has achieved substantial improvement in prompt-based zero/few-shot learning, it comes with limitations that could open avenues for future research.

First, calibration is fully unsupervised in the scenario where no labeled data is available (zero-shot downstream tasks in § 3.3). Based on empirical experimental results, we adopt the conservative *One-batch Calibration* strategy to ensure a safe and consistent performance enhancement. In the future, we aim to explore more rigorous approaches to determine optimal stopping points in this scenario.

Second, we utilize RoBERTa (encoder) models for classification tasks, as encoder models may more effectively encode task-specific patterns for discriminative tasks compared to some generative LMs (Gao et al., 2021; Li et al., 2023b), as shown in Table 16. However, the relatively small size of those Masked LMs (355M parameters for RoBERTa-large) could be the ultimate limitation to their capabilities. Given the proliferation of large-scale generative (decoder) LMs and their accomplishments in tackling more challenging tasks (Thoppilan et al., 2022; Chowdhery et al., 2023; Touvron et al., 2023), we anticipate extending our method to large decoder models and validating the applicability of our findings. Furthermore, we expect to expand the scope of tasks to include regression problems (e.g., sentiment score prediction) leveraging KL divergence to measure disparities in continuous probability distributions, aiming to address bias-related challenges across diverse scenarios.

## 8 Ethics Statement and Broader Impact

Our work is conformant to the Code of Ethics. We appropriately cite relevant methods, models, and datasets that we use. We affirm that all datasets in our experiments are public, and no private or sensitive information is incorporated in our research. Our use of datasets and pre-trained models is consistent with their intended use. For broader impacts, our method, extending beyond calibrating common token bias and association bias, might inspire prospective research in mitigating social bias and improving the fairness of pre-trained LMs.

## References

- Ankur Bapna and Orhan Firat. 2019. [Simple, scalable adaptation for neural machine translation](#). 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 1538–1548, Hong Kong, China. Association for Computational Linguistics.
- Elad Ben Zaken, Yoav Goldberg, and Shauli Ravfogel. 2022. [BitFit: Simple parameter-efficient fine-tuning for transformer-based masked language-models](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1–9, Dublin, Ireland. Association for Computational Linguistics.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Jiahao Cao, Rui Liu, Huailiang Peng, Lei Jiang, and Xu Bai. 2022. Aspect is not you need: No-aspect differential sentiment framework for aspect-based sentiment analysis. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1599–1609.
- Pengyu Cheng, Weituo Hao, Siyang Yuan, Shijing Si, and Lawrence Carin. 2021. Fairfil: Contrastive neural debiasing method for pretrained text encoders. *arXiv preprint arXiv:2103.06413*.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Emily Dinan, Angela Fan, Adina Williams, Jack Urbanek, Douwe Kiela, and Jason Weston. 2020.

673	Queens are powerful too: Mitigating gender bias in dialogue generation. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 8173–8188, Online. Association for Computational Linguistics.	729
674		730
675		731
676		732
677		733
678	William B. Dolan and Chris Brockett. 2005. <b>Automatically constructing a corpus of sentential paraphrases.</b> In <i>Proceedings of the Third International Workshop on Paraphrasing (IWP2005)</i> .	734
679		735
680		736
681		737
682	Li Dong, Furu Wei, Chuanqi Tan, Duyu Tang, Ming Zhou, and Ke Xu. 2014. <b>Adaptive recursive neural network for target-dependent Twitter sentiment classification.</b> In <i>Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 49–54, Baltimore, Maryland. Association for Computational Linguistics.	738
683		739
684		740
685		741
686		742
687		743
688		744
689		745
690	Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. <b>Making pre-trained language models better few-shot learners.</b> In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 3816–3830, Online. Association for Computational Linguistics.	746
691		747
692		748
693		749
694		750
695		751
696		752
697		753
698	Yuxian Gu, Xu Han, Zhiyuan Liu, and Minlie Huang. 2022. <b>PPT: Pre-trained prompt tuning for few-shot learning.</b> In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8410–8423, Dublin, Ireland. Association for Computational Linguistics.	754
699		755
700		756
701		757
702		758
703		759
704	Xu Guo, Boyang Li, and Han Yu. 2022. Improving the sample efficiency of prompt tuning with domain adaptation. <i>arXiv preprint arXiv:2210.02952</i> .	760
705		761
706		762
707	Zhixiong Han, Yaru Hao, Li Dong, Yutao Sun, and Furu Wei. 2023. <b>Prototypical calibration for few-shot learning of language models.</b> In <i>The Eleventh International Conference on Learning Representations</i> .	763
708		764
709		765
710		766
711		767
712	Ari Holtzman, Peter West, Vered Shwartz, Yejin Choi, and Luke Zettlemoyer. 2022. Surface form competition: Why the highest probability answer isn’t always right. <i>arXiv preprint arXiv:2104.08315</i> .	768
713		769
714		770
715		771
716	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In <i>International conference on machine learning</i> , pages 2790–2799. PMLR.	772
717		773
718		774
719		775
720		776
721		777
722	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	778
723		779
724		780
725		781
726		782
727	Po-Sen Huang, Huan Zhang, Ray Jiang, Robert Stanforth, Johannes Welbl, Jack Rae, Vishal Maini, Dani Yogatama, and Pushmeet Kohli. 2020. <b>Reducing sentiment bias in language models via counterfactual evaluation.</b> In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 65–83, Online. Association for Computational Linguistics.	783
728		784
	Yiren Jian, Chongyang Gao, and Soroush Vosoughi. 2022. <b>Contrastive learning for prompt-based few-shot language learners.</b> In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5577–5587, Seattle, United States. Association for Computational Linguistics.	785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000

785	<a href="#">Cutting down on prompts and parameters: Simple few-shot learning with language models.</a>	842
786	In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2824–2835, Dublin, Ireland. Association for Computational Linguistics.	843
787		844
788		845
789		846
790	Daniel Loureiro, Francesco Barbieri, Leonardo Neves, Luis Espinosa Anke, and Jose Camacho-collados. 2022. <a href="#">TimeLMs: Diachronic language models from Twitter.</a>	847
791	In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> , pages 251–260, Dublin, Ireland. Association for Computational Linguistics.	848
792		849
793		850
794		851
795		852
796		853
797	Jinghui Lu, Dongsheng Zhu, Weidong Han, Rui Zhao, Brian Mac Namee, and Fei Tan. 2023. <a href="#">What makes pre-trained language models better zero-shot learners?</a>	854
798	In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2288–2303, Toronto, Canada. Association for Computational Linguistics.	855
799		856
800		857
801		858
802		859
803		860
804	Nicholas Meade, Elinor Poole-Dayana, and Siva Reddy. 2022. <a href="#">An empirical survey of the effectiveness of debiasing techniques for pre-trained language models.</a>	861
805	In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1878–1898, Dublin, Ireland. Association for Computational Linguistics.	862
806		863
807		864
808		865
809		866
810		867
811	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2017. <a href="#">Pointer sentinel mixture models.</a>	868
812	In <i>International Conference on Learning Representations</i> .	869
813		870
814		871
815	Sewon Min, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. <a href="#">Noisy channel language model prompting for few-shot text classification.</a>	872
816	In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5316–5330, Dublin, Ireland. Association for Computational Linguistics.	873
817		874
818		875
819		876
820		877
821		878
822	Bo Pang and Lillian Lee. 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. <i>arXiv preprint cs/0409058</i> .	879
823		880
824		881
825	Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. 2016. <a href="#">The LAMBADA dataset: Word prediction requiring a broad discourse context.</a>	882
826	In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1525–1534, Berlin, Germany. Association for Computational Linguistics.	883
827		884
828		885
829		886
830		887
831		888
832		889
833		890
834	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. <i>Advances in neural information processing systems</i> , 32.	891
835		892
836		893
837		894
838		895
839		896
840	Maria Pontiki, Dimitris Galanis, John Pavlopoulos, Harris Papageorgiou, Ion Androutsopoulos, and Suresh	897
841	Manandhar. 2014. <a href="#">SemEval-2014 task 4: Aspect based sentiment analysis.</a>	898
	In <i>Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)</i> , pages 27–35, Dublin, Ireland. Association for Computational Linguistics.	899
		900
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	901
		902
	Julian Salazar, Davis Liang, Toan Q. Nguyen, and Katrin Kirchhoff. 2020. <a href="#">Masked language model scoring.</a>	903
	In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2699–2712, Online. Association for Computational Linguistics.	904
		905
	Timo Schick and Hinrich Schütze. 2021a. <a href="#">Exploiting cloze-questions for few-shot text classification and natural language inference.</a>	906
	In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 255–269, Online. Association for Computational Linguistics.	907
		908
	Timo Schick and Hinrich Schütze. 2021b. <a href="#">It’s not just size that matters: Small language models are also few-shot learners.</a>	909
	In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2339–2352, Online. Association for Computational Linguistics.	910
		911
	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In <i>Proceedings of the 2013 conference on empirical methods in natural language processing</i> , pages 1631–1642.	912
		913
	Irene Solaiman and Christy Dennison. 2021. Process for adapting language models to society (palms) with values-targeted datasets. <i>Advances in Neural Information Processing Systems</i> , 34:5861–5873.	914
		915
	Yi Sun, Yu Zheng, Chao Hao, and Hangping Qiu. 2022. <a href="#">NSP-BERT: A prompt-based few-shot learner through an original pre-training task — next sentence prediction.</a>	916
	In <i>Proceedings of the 29th International Conference on Computational Linguistics</i> , pages 3233–3250, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.	917
		918
	Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. 2022. <a href="#">Lamda: Language models for dialog applications.</a> <i>arXiv preprint arXiv:2201.08239</i> .	919
		920
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. <a href="#">Llama: Open and efficient foundation language models.</a> <i>arXiv preprint arXiv:2302.13971</i> .	921
		922



900 Laorens van der Maaten and Geoffrey Hinton. 2008. *Annual Meeting of the Association for Computational* 957  
901 [Visualizing data using t-sne](#). *Journal of Machine* 958  
902 *Learning Research*, 9(86):2579–2605. 4241. 959

903 Ellen M Voorhees and Dawn M Tice. 2000. Building a  
904 question answering test collection. In *Proceedings*  
905 *of the 23rd annual international ACM SIGIR confer-*  
906 *ence on Research and development in information*  
907 *retrieval*, pages 200–207.

908 Alex Wang, Amanpreet Singh, Julian Michael, Felix  
909 Hill, Omer Levy, and Samuel Bowman. 2018. [GLUE:](#)  
910 [A multi-task benchmark and analysis platform for nat-](#)  
911 [ural language understanding](#). In *Proceedings of the*  
912 *2018 EMNLP Workshop BlackboxNLP: Analyzing*  
913 *and Interpreting Neural Networks for NLP*, pages  
914 353–355, Brussels, Belgium. Association for Com-  
915 putational Linguistics.

916 Adina Williams, Nikita Nangia, and Samuel Bowman.  
917 2018. [A broad-coverage challenge corpus for sen-](#)  
918 [tence understanding through inference](#). In *Proceed-*  
919 *ings of the 2018 Conference of the North American*  
920 *Chapter of the Association for Computational Lin-*  
921 *guistics: Human Language Technologies, Volume*  
922 *1 (Long Papers)*, pages 1112–1122, New Orleans,  
923 Louisiana. Association for Computational Linguis-  
924 tics.

925 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien  
926 Chaumond, Clement Delangue, Anthony Moi, Pier-  
927 ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-  
928 icz, Joe Davison, Sam Shleifer, Patrick von Platen,  
929 Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,  
930 Teven Le Scao, Sylvain Gugger, Mariama Drame,  
931 Quentin Lhoest, and Alexander Rush. 2020. [Trans-](#)  
932 [formers: State-of-the-art natural language processing](#).  
933 In *Proceedings of the 2020 Conference on Empirical*  
934 *Methods in Natural Language Processing: System*  
935 *Demonstrations*, pages 38–45, Online. Association  
936 for Computational Linguistics.

937 Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng  
938 Huang. 2022. [NoisyTune: A little noise can help](#)  
939 [you finetune pretrained language models better](#). In  
940 *Proceedings of the 60th Annual Meeting of the As-*  
941 *sociation for Computational Linguistics (Volume 2:*  
942 *Short Papers)*, pages 680–685, Dublin, Ireland. As-  
943 sociation for Computational Linguistics.

944 Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015.  
945 Character-level convolutional networks for text classi-  
946 fication. *Advances in neural information processing*  
947 *systems*, 28.

948 Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and  
949 Sameer Singh. 2021. Calibrate before use: Improv-  
950 ing few-shot performance of language models. In  
951 *International Conference on Machine Learning*, pages  
952 12697–12706. PMLR.

953 Fan Zhou, Yuzhou Mao, Liu Yu, Yi Yang, and Ting  
954 Zhong. 2023. Causal-debias: Unifying debiasing  
955 in pretrained language models and fine-tuning via  
956 causal invariant learning. In *Proceedings of the 61st*



## A Experimental Details

**Prompts with or without demonstrations.** Table 7 shows the prompt templates and label words of each dataset we use for main experiments.

For downstream tasks, in few-shot setting, task-specific example-label pairs (i.e., demonstrations) can be incorporated in the context to enhance the LM’s comprehension. While in zero-shot setting, no labeled data is available and thereby no demonstrations.

For calibration, demonstrations are either absent from or added to null-input prompts, consistent with their exclusion from or inclusion in prompts for downstream tasks. An example of a null-input prompt without demonstration is:

`<s> An empty sentence. It is <mask>. </s>`

`<s>` and `</s>` respectively denote `<cls>` token and `<sep>` token in RoBERTa. In the other case, we incorporate demonstrations retrieved from the small training set into the null-input prompt such as:

`<s> An empty sentence. It is <mask>. </s>`

`Compellingly watchable. It is great. </s>`

`The film is strictly routine. It is terrible. </s>`

### Association-bias calibration for aspect-level task.

For aspect-level sentiment analysis, e.g., "Wonderful food but poor service. Service was <mask>.", the answer contains the aspect word "service". Because the model makes sentiment predictions for specific aspect words, the task is likely subject to *association bias* (§ 2). For association-bias calibration, the only difference is that we incorporate various aspect words in the answer format (e.g., "<aspect words> was <mask>.") when constructing null-input prompts. One can either leverage GPT-4 to generate in-domain aspect words (e.g., for restaurant reviews, the generated aspect words could be *menu*, *food*, etc.), or simply employ the aspect words in the original training dataset. In this work, we choose the latter option. Due to the variability of <aspect words> in the answer format, sorting null-meaning inputs by NSP score can yield different results. To this effect, we do not apply  $x_{\text{null}}$  selection strategy (§ 4.2) for aspect-level task, and instead keep all the generated  $x_{\text{null}}$ .

### Null-meaning inputs for One-batch Calibration.

For zero-shot downstream tasks, since only one batch of null-meaning inputs is required for calibration in our early-stopping criterion (§ 3.3), we select the  $Top-N\{P_{\text{NSP}}(x_{\text{null}}, \text{ans})\}$   $x_{\text{null}}$  from  $\mathcal{X}_{\text{null}}$ ,

where  $N$  is batch size. We prioritize these samples as our observations show that null-meaning inputs with higher  $P_{\text{NSP}}(x_{\text{null}}, \text{ans})$  exhibit higher attention scores between the null input and `<mask>`, as demonstrated in Figure 5. This indicates more effective conveyance of the "null" information to the placeholder `<mask>`, which could facilitate LM deciphering the "null" patterns of the prompts and benefit calibration.

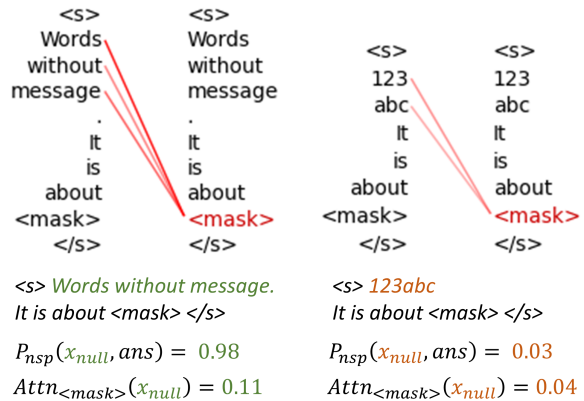


Figure 5: Visualization of attention score by the depth of color in the connecting lines. We only show the attention between `<mask>` token and null-meaning input  $x_{\text{null}}$ .  $\text{Attn}_{\text{<mask>}}(x_{\text{null}})$  is the attention score of `<mask>` on  $x_{\text{null}}$ , averaged over encoder layers and attention heads. **Left:** Higher attention score indicates enhanced pattern extraction from  $x_{\text{null}}$  which has higher  $P_{\text{NSP}}(x_{\text{null}}, \text{ans})$ .

**Hyper-parameters.** In calibration stage, we shuffle the null-input prompts and conduct gradient descent on  $\mathbf{B}_{LM}$  (or  $\mathbf{W}_{LM} + \mathbf{B}_{LM}$  as comparative experiment) with 5 different seeds to account for calibration variance. There are two main hyper-parameters for calibration: (1)  $x_{\text{null}}$  batch size  $N$ ; (2) calibration learning rate  $lr_{\text{calib}}$ . We conduct grid search on  $N = \{8, 16, 32\}$  and  $lr_{\text{calib}}$  from  $1e - 6$  to  $1e - 3$ , and obtain the best settings:  $N = 32$  and  $lr_{\text{calib}}$  as shown in Table 6.

Calibrated LMs are applied in downstream tasks with prompt-learning methods. We use the same hyper-parameters as Gao et al. (2021) for prompt learning. We evaluate on each task’s original test set, except for AGNews and DBpedia, where we randomly sample 2000 test examples.

We use PyTorch (Paszke et al., 2019) and public HuggingFace Transformers library (Wolf et al., 2020), and conduct all the experiments with one NVIDIA V100 GPU in Google Colab.

	Calibration ( $lr_{calib}$ )		Prompt FT (downstream)
	$W_{LM} + B_{LM}$	$B_{LM}$	
No demo	$1e-5$	$1e-3$	$1e-5$
With demo	$1e-6$	$1e-4$	$1e-5$

Table 6: Optimal learning rates for calibration and downstream prompt-based fine-tuning (Prompt FT). With/No demo denotes adding/not adding demonstrations in prompts.

---

**Algorithm 1** Null-input prompting for calibration

---

**Inputs:**

Downstream task: *zero\_shot* or *few\_shot*

Null-input prompts:  $\{N_{prompt}\}$

(Val. data in Calibration:  $\mathcal{D}_{val}^{calib} \leftarrow \mathcal{D}_{train}^{downstream}$ )

▷ Only when downstream task is *few\_shot*.

▷ Downstream training dataset  $\mathcal{D}_{train}^{downstream}$  constitutes  $K$  samples per class.

**Output:**

$LM_{calib}^{one\_batch}$  for *zero\_shot*

$LM_{calib}^{one\_batch}$  &  $LM_{calib}^{val}$  for *few\_shot*

- 1: **for**  $batch$  in  $\{N_{prompt}\}$  **do**
  - 2:      $P = \mathcal{LM}(batch)$  ▷ Null input prompting
  - 3:      $\mathcal{L} = D_{KL}(U || P)$  ▷ Unif. distribution  $U$
  - 4:      $B_{LM} \leftarrow B_{LM} - \alpha \cdot \frac{\partial \mathcal{L}}{\partial B_{LM}}$  ▷ Freeze  $W_{LM}$
  - 5:     **if** *first batch* **then**
  - 6:         Save  $LM_{calib}^{one\_batch}$
  - 7:     **end if**
  - 8:     **if** downstream is *zero\_shot* **then break**
  - 9:     **end if**
  - 10:    **if** better  $Compute\_Metric(\mathcal{D}_{val}^{calib})$  **then**
  - 11:        Save  $LM_{calib}^{val}$
  - 12:    **end if**
  - 13: **end for**
- 

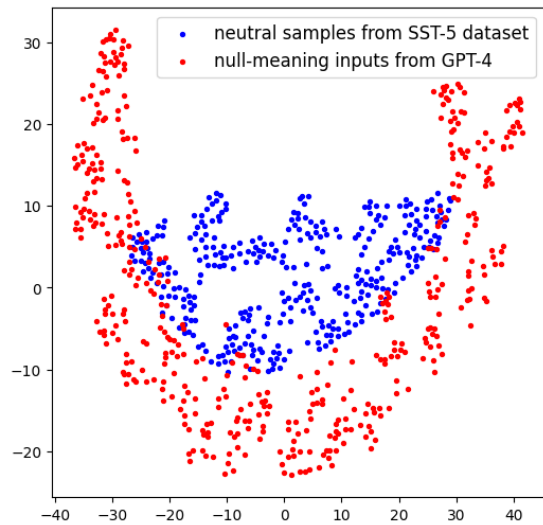


Figure 6: t-SNE visualization of output representations for null-meaning inputs generated from GPT-4 (red) compared to neutral samples from SST-5 dataset (blue). We utilize the pre-trained sentiment analysis model (Loureiro et al., 2022) to obtain the embeddings. The different distributions validate that null information is not equivalent to neutral sentiment.

Dataset	Task Type	Prompt Template	Label Words
AGNews	News topic classification	{Sentence} It is about <mask>.	World / Sports / Business / Technology
DBPedia <sup>†</sup>	Ontology classification	{Sentence} It is about <mask>.	Company / Artist / Building / Nature
TREC	Question classification	{Sentence} It is about <mask>.	Number / Location / Person / Description / Entity / Expression
Subj	Subjectivity classification	{Sentence} This is <mask>.	objective / subjective
SST-5	Movie sentiment analysis	{Sentence} The movie was <mask>.	terrible / bad / okay / good / great
Laptop	Aspect level sentiment analysis	{Sentence} {Aspect words} was <mask>.	terrible / okay / great
Restaurant	Aspect level sentiment analysis	{Sentence} {Aspect words} was <mask>.	terrible / okay / great
Twitter	Aspect level sentiment analysis	{Sentence} {Aspect words} was <mask>.	terrible / okay / great

Table 7: Prompt templates and label words of the eight datasets in our experiments for main results. For DBPedia<sup>†</sup>, we use four classes out of the total fourteen classes.

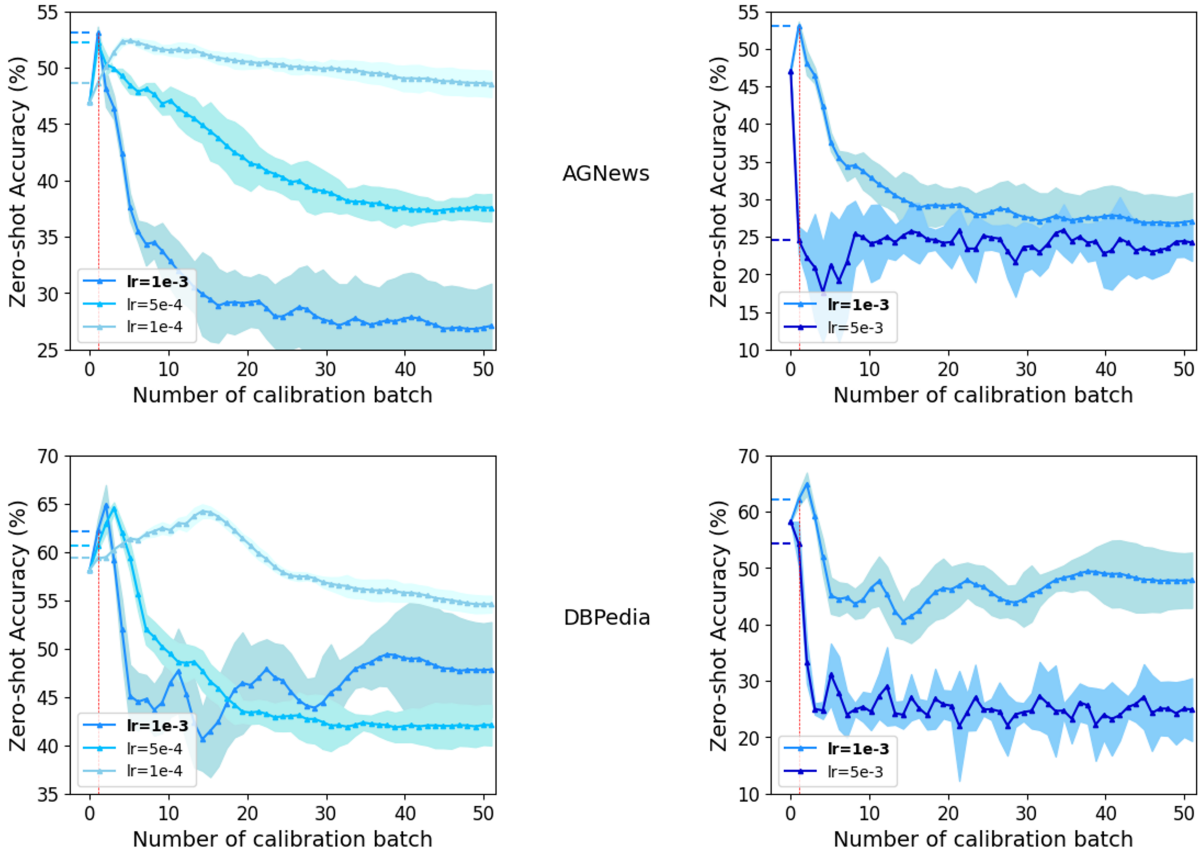


Figure 7: Empirical experiments show the impact of calibration on zero-shot learning performance across *different calibration learning rates*  $lr_{calib}$ , with a fixed batch size of 32. Only  $\mathbf{B}_{LM}$  is updated in calibration. We identify the optimal  $lr_{calib} = 1e - 3$  across all datasets and illustrate with AGNews dataset (top two figures) and DBPedia dataset (bottom two figures). A smaller learning rate (left figures) consistently yields less performance improvement, considering both peak accuracy and accuracy after the first calibration batch (the intersections of the curves and red vertical line). A larger learning rate (right figures) consistently degrades performance.

## B Additional Results

### B.1 Performance Comparison with NSP-BERT, Perplection and NoisyTune

We additionally choose NSP-BERT (Sun et al., 2022) and Perplection (Lu et al., 2023) as *in-context learning* comparison baselines and NoisyTune (Wu et al., 2022) as *prompt-based fine-tuning* comparison baseline. NSP-BERT constructs potential answers using each label word and predict Next Sentence Prediction (NSP) probability between the input and each answer. Perplection proposes perplexity-based selection method for zero-shot prompt learning. NoisyTune demonstrates that adding noise to pre-trained LMs benefits fine-tuning on downstream tasks. We re-implement their methods with the same settings as ours for fair comparisons. As shown in Table 8 and Table 9, our method achieves superior results in almost all datasets.

Furthermore, our method consistently outperforms NoisyTune, demonstrating that the gains in prompt-based fine-tuning with our method are not solely a result of perturbing LM parameters. This confirms the efficacy of intrinsic bias calibration in enhancing LM performance.

	Zero-shot in-context learning		
	NSP-BERT	Perplection	IntrCal
AGNews	52.4	49.3	<b>54.5</b>
DBPedia	58.4	59.6	<b>61.8</b>
TREC	<b>32.4</b>	30.8	31.1
Subj	60.3	59.9	<b>62.7</b>
SST-5	30.2	31.0	<b>37.5</b>
Laptop	57.3	58.2	<b>59.6</b>
Restaurant	50.4	66.5	<b>72.8</b>
Twitter	35.3	31.5	<b>51.7</b>
Average	47.1	48.4	<b>54.0</b>

Table 8: Comparison of NSP-BERT (Sun et al., 2022), Perplection (Lu et al., 2023) and IntrCal (ours) in zero-shot in-context learning.

### B.2 Other Experiments

We briefly summarize the contents of each table and figure below that presents other additional results.

Table 10 contains results for performance using RoBERTa-base model.

Table 11 contains results for performance of  $K = \{2, 4, 8\}$  few-shot learning.

	Prompt FT no demo		Prompt FT with demo	
	NoisyTune	IntrCal	NoisyTune	IntrCal
AGNews	89.0 <sub>1.8</sub>	<b>89.0</b> <sub>0.8</sub>	88.4 <sub>1.5</sub>	<b>89.3</b> <sub>0.9</sub>
DBPedia	98.0 <sub>0.8</sub>	<b>99.0</b> <sub>0.1</sub>	98.6 <sub>0.9</sub>	<b>98.9</b> <sub>0.3</sub>
TREC	86.2 <sub>4.3</sub>	<b>89.3</b> <sub>4.5</sub>	87.2 <sub>4.6</sub>	<b>89.7</b> <sub>1.0</sub>
Subj	93.0 <sub>1.2</sub>	<b>93.2</b> <sub>1.2</sub>	92.9 <sub>1.2</sub>	<b>94.3</b> <sub>0.2</sub>
SST-5	49.4 <sub>1.1</sub>	<b>49.9</b> <sub>2.7</sub>	47.5 <sub>3.5</sub>	<b>50.0</b> <sub>1.7</sub>
Laptop	73.8 <sub>3.2</sub>	<b>74.9</b> <sub>2.9</sub>	75.5 <sub>3.2</sub>	<b>78.7</b> <sub>1.4</sub>
Restaurant	79.9 <sub>2.7</sub>	<b>82.0</b> <sub>0.9</sub>	78.3 <sub>2.6</sub>	<b>79.8</b> <sub>4.5</sub>
Twitter	51.8 <sub>5.8</sub>	<b>57.0</b> <sub>4.2</sub>	59.0 <sub>1.9</sub>	<b>59.3</b> <sub>2.3</sub>
Average	77.6	<b>79.3</b>	78.4	<b>80.0</b>

Table 9: Comparison between NoisyTune (Wu et al., 2022) and IntrCal (ours) in prompt-based fine-tuning.

Figure 8 contains results for performance using different prompt templates (Table 13).

Table 12 contains results for pseudo-perplexity comparisons between updating entire LM and only updating bias parameters in calibration.

Table 14 contains results for performance comparisons between updating entire LM and only updating bias parameters in calibration.

Table 15 contains results for performance of sentence-pair datasets.

Table 16 contains results for performance comparisons between Llama-2 and RoBERTa.

Table 17 contains results for variance of probability distribution across labels before and after calibration.



	In-context lrn no demo			In-context lrn with demo			Prompt FT no demo			Prompt FT with demo		
	NoCal	OutCal	IntrCal	NoCal	OutCal	IntrCal	NoCal	OutCal	IntrCal	NoCal	OutCal	IntrCal
AGNews	37.8 <sub>0.0</sub>	36.2 <sub>4.6</sub>	<b>49.0</b> <sub>0.9</sub>	68.4 <sub>0.4</sub>	69.7 <sub>4.3</sub>	<b>73.7</b> <sub>0.3</sub>	88.2 <sub>0.3</sub>	87.8 <sub>0.6</sub>	<b>88.9</b> <sub>1.0</sub>	86.7 <sub>0.1</sub>	74.2 <sub>4.1</sub>	<b>87.2</b> <sub>0.1</sub>
DBPedia	<b>57.2</b> <sub>0.0</sub>	50.5 <sub>7.1</sub>	54.9 <sub>0.1</sub>	56.5 <sub>3.4</sub>	78.7 <sub>4.4</sub>	<b>83.9</b> <sub>0.4</sub>	95.2 <sub>2.1</sub>	93.5 <sub>5.0</sub>	<b>99.0</b> <sub>0.4</sub>	97.8 <sub>0.9</sub>	96.7 <sub>0.8</sub>	<b>98.6</b> <sub>0.1</sub>
TREC	28.2 <sub>0.0</sub>	25.4 <sub>4.4</sub>	<b>30.2</b> <sub>0.1</sub>	41.2 <sub>0.3</sub>	39.9 <sub>3.8</sub>	<b>42.5</b> <sub>1.0</sub>	82.5 <sub>10.9</sub>	70.3 <sub>2.3</sub>	<b>86.4</b> <sub>6.5</sub>	85.7 <sub>1.8</sub>	80.6 <sub>5.0</sub>	<b>91.2</b> <sub>0.6</sub>
Subj	53.6 <sub>0.0</sub>	63.6 <sub>1.9</sub>	<b>66.4</b> <sub>1.8</sub>	50.8 <sub>0.2</sub>	67.0 <sub>1.7</sub>	<b>69.6</b> <sub>0.4</sub>	<b>92.5</b> <sub>1.3</sub>	91.1 <sub>0.4</sub>	91.9 <sub>1.7</sub>	90.4 <sub>2.1</sub>	92.0 <sub>0.2</sub>	<b>92.3</b> <sub>0.1</sub>
SST-5	31.9 <sub>0.0</sub>	30.8 <sub>3.4</sub>	<b>32.2</b> <sub>0.2</sub>	25.3 <sub>3.4</sub>	28.6 <sub>3.4</sub>	<b>29.8</b> <sub>1.7</sub>	45.9 <sub>3.3</sub>	42.9 <sub>2.3</sub>	<b>48.1</b> <sub>1.8</sub>	44.3 <sub>5.2</sub>	40.7 <sub>2.5</sub>	<b>45.8</b> <sub>2.6</sub>
Laptop	56.1 <sub>0.0</sub>	56.7 <sub>3.8</sub>	<b>60.0</b> <sub>0.1</sub>	49.2 <sub>0.9</sub>	61.5 <sub>2.8</sub>	<b>64.0</b> <sub>0.6</sub>	75.8 <sub>3.4</sub>	73.0 <sub>1.3</sub>	<b>76.3</b> <sub>1.8</sub>	74.8 <sub>0.1</sub>	76.0 <sub>0.6</sub>	<b>76.3</b> <sub>0.5</sub>
Restaurant	69.8 <sub>0.0</sub>	<b>72.0</b> <sub>2.9</sub>	69.5 <sub>0.5</sub>	67.6 <sub>0.7</sub>	70.5 <sub>2.4</sub>	<b>73.2</b> <sub>0.7</sub>	75.5 <sub>6.6</sub>	<b>77.3</b> <sub>3.4</sub>	77.2 <sub>1.1</sub>	74.8 <sub>3.3</sub>	75.2 <sub>0.7</sub>	<b>76.1</b> <sub>3.9</sub>
Twitter	22.0 <sub>0.0</sub>	48.6 <sub>5.1</sub>	<b>52.3</b> <sub>0.6</sub>	17.6 <sub>0.4</sub>	41.8 <sub>5.4</sub>	<b>48.4</b> <sub>0.5</sub>	54.5 <sub>1.1</sub>	47.7 <sub>3.8</sub>	<b>57.9</b> <sub>1.3</sub>	50.6 <sub>4.6</sub>	51.8 <sub>2.1</sub>	<b>56.0</b> <sub>4.9</sub>
Average	44.6	48.0	<b>51.8</b>	47.1	57.2	<b>60.6</b>	76.3	73.0	<b>78.2</b>	75.6	73.4	<b>77.9</b>

Table 10: Result comparisons among NoCal (LM-BFF Gao et al., 2021; no calibration), OutCal (output calibration) and IntrCal (ours; intrinsic-bias calibrated LM) using RoBERTa-base. We report the mean and standard deviation of performance in 8 classification datasets with 4 prompt-learning methods.

		In-context lrn with demo		Prompt FT no demo		Prompt FT with demo	
		NoCal	IntrCal	NoCal	IntrCal	NoCal	IntrCal
2-shot	AGNews	70.4 <sub>6.7</sub>	<b>76.3</b> <sub>3.6</sub>	76.4 <sub>5.4</sub>	<b>80.2</b> <sub>8.0</sub>	78.2 <sub>1.3</sub>	<b>83.2</b> <sub>1.1</sub>
	DBPedia	92.9 <sub>0.9</sub>	<b>94.0</b> <sub>1.0</sub>	97.0 <sub>1.6</sub>	<b>98.4</b> <sub>0.9</sub>	97.4 <sub>1.0</sub>	<b>97.8</b> <sub>1.1</sub>
	TREC	49.8 <sub>4.2</sub>	<b>50.5</b> <sub>4.0</sub>	49.1 <sub>22.6</sub>	<b>60.3</b> <sub>9.6</sub>	65.2 <sub>9.3</sub>	<b>66.1</b> <sub>9.3</sub>
	Subj	49.4 <sub>1.1</sub>	<b>56.2</b> <sub>3.9</sub>	66.4 <sub>5.4</sub>	<b>82.2</b> <sub>5.9</sub>	72.3 <sub>13.9</sub>	<b>81.5</b> <sub>13.2</sub>
4-shot	AGNews	75.7 <sub>3.9</sub>	<b>80.3</b> <sub>1.7</sub>	85.4 <sub>2.7</sub>	<b>87.3</b> <sub>1.3</sub>	76.7 <sub>13.1</sub>	<b>85.9</b> <sub>1.9</sub>
	DBPedia	93.0 <sub>0.4</sub>	<b>93.9</b> <sub>0.4</sub>	97.2 <sub>0.8</sub>	<b>97.9</b> <sub>1.1</sub>	96.4 <sub>1.5</sub>	<b>98.6</b> <sub>0.6</sub>
	TREC	51.9 <sub>2.6</sub>	<b>53.2</b> <sub>2.5</sub>	64.5 <sub>7.1</sub>	<b>67.6</b> <sub>6.7</sub>	73.6 <sub>8.5</sub>	<b>78.2</b> <sub>9.7</sub>
	Subj	48.8 <sub>2.2</sub>	<b>59.4</b> <sub>3.1</sub>	81.4 <sub>3.9</sub>	<b>88.5</b> <sub>3.2</sub>	78.9 <sub>9.3</sub>	<b>83.6</b> <sub>7.8</sub>
8-shot	AGNews	79.6 <sub>1.0</sub>	<b>82.4</b> <sub>1.6</sub>	86.9 <sub>1.9</sub>	<b>88.1</b> <sub>0.4</sub>	85.5 <sub>1.7</sub>	<b>88.0</b> <sub>1.4</sub>
	DBPedia	92.9 <sub>0.8</sub>	<b>94.2</b> <sub>0.2</sub>	97.3 <sub>1.2</sub>	<b>98.8</b> <sub>0.5</sub>	98.2 <sub>0.8</sub>	<b>98.6</b> <sub>0.2</sub>
	TREC	47.9 <sub>2.2</sub>	<b>48.7</b> <sub>2.0</sub>	71.6 <sub>4.9</sub>	<b>72.2</b> <sub>5.1</sub>	75.4 <sub>6.2</sub>	<b>81.7</b> <sub>5.6</sub>
	Subj	48.4 <sub>1.0</sub>	<b>60.5</b> <sub>4.8</sub>	91.9 <sub>1.3</sub>	<b>92.7</b> <sub>0.8</sub>	88.9 <sub>5.3</sub>	<b>92.1</b> <sub>2.2</sub>

Table 11: Few-shot learning with different number of training samples ( $K = \{2, 4, 8\}$ ) using RoBERTa-large. IntrCal (ours; intrinsic-bias calibrated LM) consistently outperforms NoCal (no calibration).

Model	Datasets					
	WikiText-2		WikiText-103		LAMBADA	
Original RoBERTa	6.189		7.008		24.52	
+ CALIBRATION	$W_{LM} + B_{LM}$	$B_{LM}$	$W_{LM} + B_{LM}$	$B_{LM}$	$W_{LM} + B_{LM}$	$B_{LM}$
for_AGNews	$\uparrow$ 0.105 6.294	$\uparrow$ 0.017 6.206	$\uparrow$ 0.059 7.067	$\uparrow$ 0.029 7.037	$\uparrow$ 0.58 25.10	$\uparrow$ 0.02 24.54
for_DBPedia	$\uparrow$ 0.101 6.290	$\uparrow$ 0.008 6.197	$\uparrow$ 0.092 7.100	$\uparrow$ 0.002 7.010	$\uparrow$ 0.76 25.28	$\downarrow$ 0.22 24.30
for_TREC	$\uparrow$ 0.049 6.238	$\downarrow$ 0.027 6.162	$\uparrow$ 0.040 7.048	$\downarrow$ 0.042 6.966	$\uparrow$ 0.57 25.09	$\downarrow$ 0.27 24.25
for_Subj	$\uparrow$ 0.081 6.270	$\downarrow$ 0.021 6.168	$\uparrow$ 0.116 7.124	$\downarrow$ 0.030 6.978	$\uparrow$ 0.70 25.22	$\uparrow$ 0.08 24.60
for_SST-5	$\downarrow$ 0.018 6.171	$\downarrow$ 0.031 6.158	$\uparrow$ 0.143 7.151	$\downarrow$ 0.039 6.969	$\uparrow$ 0.65 25.17	$\downarrow$ 0.18 24.34
for_Laptop	$\uparrow$ 0.133 6.322	$\uparrow$ 0.011 6.200	$\uparrow$ 0.075 7.083	$\uparrow$ 0.002 7.010	$\uparrow$ 0.56 25.08	$\downarrow$ 0.01 24.51
for_Restaurant	$\uparrow$ 0.102 6.291	$\uparrow$ 0.055 6.244	$\uparrow$ 0.071 7.079	$\uparrow$ 0.074 7.082	$\uparrow$ 0.64 25.16	$\uparrow$ 0.13 24.65
for_Twitter	$\uparrow$ 0.204 6.393	$\downarrow$ 0.029 6.160	$\uparrow$ 0.096 7.104	$\downarrow$ 0.037 6.971	$\uparrow$ 0.39 24.91	$\uparrow$ 0.05 24.57

Table 12: Pseudo-perplexities of original RoBERTa and task-specific calibrated RoBERTa on WikiText-2, WikiText-103 and LAMBADA. We use 2000 test samples of each dataset. An increase in values (highlighted in red) indicates a reduction in language modeling abilities after calibration.  $W_{LM} + B_{LM}$  updates entire LM in calibration while  $B_{LM}$  only updates bias parameters.

Task	Prompt Templates
AGNews	{Sentence} It is about <mask>.
	{Sentence} This is about <mask>.
	{Sentence} This is on <mask>.
	{Sentence} It pertains to <mask>.
	{Sentence} In relation to <mask>.
TREC	{Sentence} It is about <mask>.
	{Sentence} Concerning <mask>.
	{Sentence} This is about <mask>.
	{Sentence} In relation to <mask>.
	{Sentence} This is on <mask>.

Table 13: The five different prompt templates used in Figure 8.

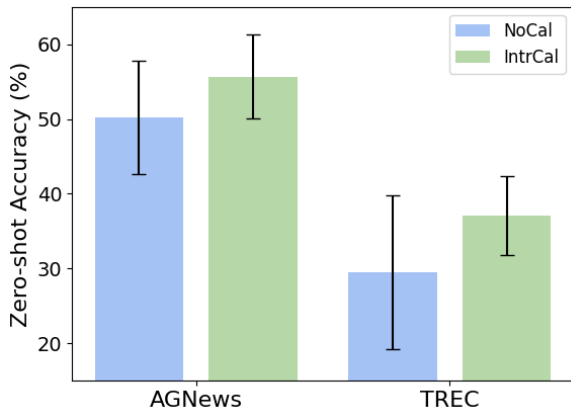


Figure 8: Performance comparison averaged on using five different prompt templates with RoBERTa-large. IntrCal (ours; intrinsic-bias calibrated LM) demonstrates significantly improved accuracy with lower variance compared to NoCal (no calibration).

	In-context lrn with demo		Prompt FT with demo	
	$W_{LM} + B_{LM}$	$B_{LM}$	$W_{LM} + B_{LM}$	$B_{LM}$
AGNews	82.0 <sub>0.8</sub>	<b>82.4</b> <sub>0.9</sub>	<b>89.3</b> <sub>0.6</sub>	89.3 <sub>0.9</sub>
DBPedia	<b>95.1</b> <sub>0.7</sub>	94.8 <sub>0.7</sub>	<b>99.0</b> <sub>0.1</sub>	98.9 <sub>0.3</sub>
TREC	<b>49.1</b> <sub>2.6</sub>	48.6 <sub>2.2</sub>	88.9 <sub>2.3</sub>	<b>89.7</b> <sub>1.0</sub>
Subj	<b>65.6</b> <sub>0.4</sub>	63.5 <sub>2.3</sub>	93.9 <sub>1.6</sub>	<b>94.3</b> <sub>0.2</sub>
SST-5	<b>37.1</b> <sub>1.0</sub>	36.6 <sub>1.0</sub>	<b>51.3</b> <sub>1.7</sub>	50.0 <sub>1.7</sub>
Laptop	65.8 <sub>0.3</sub>	<b>67.4</b> <sub>1.7</sub>	77.7 <sub>0.8</sub>	<b>78.7</b> <sub>1.4</sub>
Restaurant	72.7 <sub>1.2</sub>	<b>74.0</b> <sub>1.0</sub>	<b>81.4</b> <sub>3.4</sub>	79.8 <sub>4.5</sub>
Twitter	45.8 <sub>2.7</sub>	<b>49.4</b> <sub>2.7</sub>	<b>60.4</b> <sub>1.7</sub>	59.3 <sub>2.3</sub>
Average	64.2	<b>64.6</b>	<b>80.2</b>	80.0

Table 14: Performance comparisons between differently calibrated LMs using RoBERTa-large.  $W_{LM} + B_{LM}$  updates entire LM in calibration while  $B_{LM}$  only updates bias parameters. This table (prompt learning *with* demonstrations) is the supplement to § 5.3 Table 4 (prompt learning *without* demonstrations).

	In-context lrn no demo		Prompt FT no demo	
	NoCal	IntrCal	NoCal	IntrCal
MNLI	32.7 <sub>0.0</sub>	<b>37.7</b> <sub>0.7</sub>	67.9 <sub>2.1</sub>	<b>68.6</b> <sub>1.9</sub>
SNLI	33.6 <sub>0.0</sub>	<b>36.7</b> <sub>0.9</sub>	77.4 <sub>2.8</sub>	<b>78.5</b> <sub>2.3</sub>
MRPC	51.1 <sub>0.0</sub>	<b>53.6</b> <sub>0.2</sub>	73.6 <sub>4.3</sub>	<b>74.9</b> <sub>1.4</sub>
QQP	50.8 <sub>0.0</sub>	<b>54.6</b> <sub>0.2</sub>	65.2 <sub>3.5</sub>	<b>66.2</b> <sub>3.3</sub>

Table 15: Benchmark on sentence-pair datasets, MNLI (Williams et al., 2018), SNLI (Bowman et al., 2015), MRPC (Dolan and Brockett, 2005), QQP (Wang et al., 2018). NoCal denotes no-calibration (baseline) and IntrCal denotes our method. Our method demonstrates effectiveness on sentence-pair datasets. The overall low performance of in-context learning can be attributed to two main factors: (1) RoBERTa’s inherent limited capabilities when using in-context learning for the more difficult tasks, which is significantly improved with prompt-based fine-tuning. (2) The misalignment between these sentence-pair datasets and the use of single-sentence null inputs for calibration, which could impact the effectiveness of calibration.

	Llama-2 (7B)	RoBERTa-large (355M)
AGNews	44.1	<b>47.0</b>
DBPedia	47.1	<b>58.2</b>
TREC	<b>42.0</b>	24.0
Subj	49.8	<b>50.8</b>
SST-5	29.3	<b>31.5</b>
Laptop	48.5	<b>54.6</b>
Restaurant	65.4	<b>68.6</b>
Twitter	<b>25.5</b>	19.7

Table 16: Comparison between Llama-2 (7B parameters) (Touvron et al., 2023) and RoBERTa-large (355M parameters) on zero-shot in-context learning performance for classification tasks. Llama-2 does not consistently outperform RoBERTa in these tasks.

	AGNews	DBPedia	TREC	Subj	SST-5
Orig. LM	0.033	0.130	0.025	0.195	0.011
Calib. LM	0.022	0.025	0.011	0.112	0.011

Table 17: We calculate the **variance** of probability distribution across labels conditioned on null-meaning inputs, i.e.,  $Var(\bar{P}_{\mathcal{X}_{\text{null}}}(\mathcal{Y}))$ , before and after calibration. A smaller variance indicates that a distribution is closer to uniform distribution. Orig. LM denotes original LM, and Calib. LM denotes the LM after *One-batch Calibration* (§ 3.3). The decreasing variance in each task after calibration demonstrates that our method promotes the establishment of equitable LMs.