

# NOISYICL: A Little Noise in Model Parameters Calibrates In-context Learning

Anonymous ACL submission

## Abstract

In-Context Learning (ICL) is suffering from unsatisfactory performance and under-calibration due to high prior bias and unfaithful confidence. Some previous works fine-tuned language models for better ICL performance with enormous datasets and computing costs. In this paper, we propose NOISYICL, simply perturbing the model parameters by random noises to strive for better performance and calibration. Our experiments on two models and 12 downstream datasets show that NOISYICL can help ICL produce more accurate predictions. Our further analysis indicates that NOISYICL enables the model to provide more fair predictions, and also with more faithful confidence. Therefore, we believe that NOISYICL is an effective calibration of ICL. Our experimental code is uploaded to Github<sup>1</sup>.

## 1 Introduction

Language Models (LMs) have demonstrated the ability of In-Context Learning (ICL), where LMs learn tasks from few-shot input-label demonstrations in the form of natural language, without explicit parameter updates (Dong et al., 2022).

Nevertheless, ICL still underperforms the end-to-end models (Mosbach et al., 2023). A recent study shows that vanilla LMs are biased towards the knowledge acquired during pre-training, which is deemed harmful to ICL (Fei et al., 2023). There has been some effort in fine-tuning or calibrating LMs towards ICL tasks (Min et al., 2022; Zhao et al., 2021; Wei et al., 2021; Fei et al., 2023; Lu et al., 2022), which focus on reducing the gap between the knowledge gained from pre-training and ICL tasks, producing significant improvements in the ICL performance. However, fine-tuning these enormous LMs on additional data incurs a considerably high computational cost.

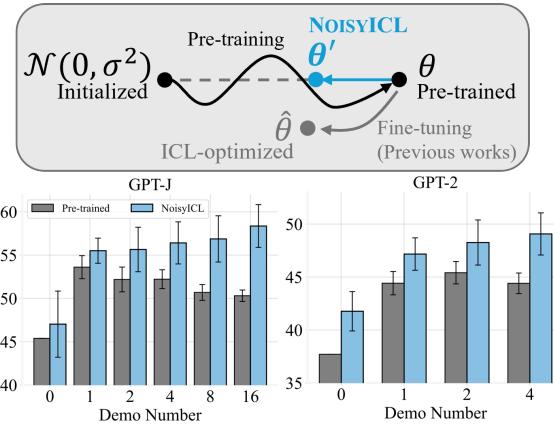


Figure 1: **Upper:** A sketch of NOISYICL: Unlike previous works which fine-tuned LMs towards ICL tasks, we perturb LMs by random noise sampled from the normal distribution  $\mathcal{N}(0, \sigma^2)$  with intensity  $\lambda$ , then perform ICL. **Lower:** The average accuracy of ICL with and without NOISYICL w.r.t. the number of demos.

Inspired by Wu et al. (2022), which shows the benefits of noise for fine-tuning end-to-end LMs, we hypothesize that introducing noise to the model parameters of pre-trained LMs can fit LMs towards ICL with a lower computational cost. In this paper, we propose NOISYICL, which simply adds noise to model parameters and then performs ICL on the noised models, as shown in Fig. 1, to investigate the benefit of noise for ICL.

As shown in Fig. 1 and Table 1, our experiments on two models and 12 classification datasets show that adding appropriate noise into model parameters improves the performance of ICL by around 10% with obviously lower computational cost.

To investigate the reason for such performance improvement, our further experiments hypothesize that NOISYICL acts as a model calibration. In detail, we find that: **1.** NOISYICL neutralizes prediction bias among general tokens and labels, and **2.** NOISYICL leads the model to produce more faithful confidence.

<sup>1</sup>Not available during anonymous review.

Table 1: Accuracy and Macro-F1 results ( $mean_{std}$ ,  $k = 4$ ). A better result is in green.  $\lambda$ : The optimal intensity of noise searched on the validation set, **Acc.**: Accuracy, **MF1**: Macro-F1, **ECE<sub>1</sub>**: the Expected Calibrated Error described in §3.3; **w/o**: Not using NOISYICL, **w/**: Using NOISYICL; Datasets abbreviation described in Appendix. **A**; **Val.**: the results on the validation sets, **Test**: the results on the testing sets.

Dataset		PS	HS	SE'14R	SE'14L	RTE	MRPC	Ethos	FP	SST2	TEE	TES	TEH	Mean
GPT-J	$\lambda$	0.2	0.2	0.1	0.1	0.012	0.2	0.004	0.08	0.004	0.01	0.08	0.002	—
	Acc. (%)	Val. w/o	37.46 <sub>1.36</sub>	71.97 <sub>1.19</sub>	36.48 <sub>0.84</sub>	34.82 <sub>0.21</sub>	49.88 <sub>1.03</sub>	43.18 <sub>1.48</sub>	53.44 <sub>0.56</sub>	62.36 <sub>0.17</sub>	81.48 <sub>0.97</sub>	45.68 <sub>0.48</sub>	49.32 <sub>1.07</sub>	50.04 <sub>1.68</sub>
		Val. w/	59.82 <sub>5.82</sub>	82.95 <sub>1.36</sub>	51.54 <sub>7.44</sub>	45.14 <sub>4.06</sub>	51.00 <sub>1.35</sub>	59.82 <sub>6.15</sub>	56.00 <sub>0.89</sub>	62.42 <sub>0.43</sub>	82.52 <sub>1.03</sub>	46.35 <sub>0.53</sub>	50.90 <sub>1.32</sub>	50.35 <sub>1.06</sub>
	Test (%)	w/o	33.09 <sub>1.76</sub>	81.30 <sub>0.79</sub>	33.67 <sub>1.31</sub>	34.26 <sub>1.33</sub>	46.52 <sub>1.11</sub>	50.97 <sub>1.43</sub>	57.74 <sub>1.38</sub>	61.87 <sub>0.16</sub>	79.15 <sub>1.17</sub>	44.72 <sub>0.62</sub>	50.20 <sub>0.85</sub>	53.16 <sub>1.07</sub>
		w/	55.90 <sub>3.43</sub>	87.16 <sub>5.29</sub>	46.79 <sub>4.47</sub>	41.77 <sub>3.80</sub>	47.64 <sub>1.27</sub>	53.45 <sub>2.23</sub>	57.07 <sub>1.53</sub>	61.90 <sub>0.30</sub>	78.82 <sub>0.75</sub>	45.17 <sub>0.57</sub>	48.11 <sub>1.38</sub>	53.12 <sub>1.10</sub>
	MF1 (%)	Val. w/o	24.11 <sub>0.87</sub>	26.75 <sub>1.38</sub>	32.53 <sub>1.46</sub>	28.62 <sub>0.68</sub>	47.37 <sub>1.14</sub>	43.18 <sub>1.48</sub>	53.24 <sub>0.69</sub>	25.91 <sub>0.26</sub>	81.24 <sub>1.04</sub>	25.65 <sub>0.54</sub>	32.87 <sub>1.81</sub>	49.01 <sub>1.81</sub>
GPT-2		Val. w/	20.80 <sub>1.13</sub>	24.22 <sub>0.63</sub>	46.88 <sub>6.70</sub>	43.25 <sub>4.12</sub>	48.61 <sub>1.38</sub>	48.38 <sub>2.69</sub>	55.89 <sub>0.86</sub>	27.81 <sub>0.74</sub>	82.33 <sub>1.06</sub>	26.31 <sub>0.80</sub>	38.19 <sub>5.33</sub>	49.37 <sub>0.92</sub>
	Test (%)	w/o	21.86 <sub>0.86</sub>	26.61 <sub>1.61</sub>	30.71 <sub>1.47</sub>	33.17 <sub>1.52</sub>	44.94 <sub>1.30</sub>	49.56 <sub>1.48</sub>	57.24 <sub>1.45</sub>	26.59 <sub>0.52</sub>	77.44 <sub>1.36</sub>	23.57 <sub>0.86</sub>	34.23 <sub>1.25</sub>	53.16 <sub>1.07</sub>
		w/	22.87 <sub>1.57</sub>	24.22 <sub>0.89</sub>	43.88 <sub>3.90</sub>	42.17 <sub>4.22</sub>	45.93 <sub>1.31</sub>	46.51 <sub>1.56</sub>	56.46 <sub>1.54</sub>	27.24 <sub>1.22</sub>	77.11 <sub>0.84</sub>	23.78 <sub>0.84</sub>	32.68 <sub>2.17</sub>	53.11 <sub>1.10</sub>
	$ECE_1$ (%)	Val. w/o	28.70 <sub>1.78</sub>	14.39 <sub>1.13</sub>	31.20 <sub>0.69</sub>	38.56 <sub>0.36</sub>	29.35 <sub>0.95</sub>	27.79 <sub>1.50</sub>	23.20 <sub>0.62</sub>	23.90 <sub>0.17</sub>	11.45 <sub>0.93</sub>	42.79 <sub>0.36</sub>	17.57 <sub>0.81</sub>	28.45 <sub>1.23</sub>
		Val. w/	12.77 <sub>2.94</sub>	9.87 <sub>2.56</sub>	13.96 <sub>7.26</sub>	21.34 <sub>6.75</sub>	28.11 <sub>0.95</sub>	8.79 <sub>3.41</sub>	20.33 <sub>0.99</sub>	21.43 <sub>3.83</sub>	12.86 <sub>1.06</sub>	42.03 <sub>0.88</sub>	10.28 <sub>5.37</sub>	28.18 <sub>1.37</sub>
	Test (%)	w/o	30.33 <sub>1.42</sub>	9.41 <sub>0.78</sub>	34.45 <sub>1.39</sub>	36.78 <sub>1.56</sub>	32.42 <sub>1.10</sub>	28.29 <sub>1.47</sub>	18.97 <sub>1.45</sub>	23.70 <sub>0.49</sub>	10.75 <sub>1.02</sub>	44.09 <sub>0.79</sub>	16.56 <sub>0.56</sub>	24.50 <sub>0.99</sub>
		w/	9.65 <sub>2.83</sub>	5.97 <sub>3.45</sub>	17.53 <sub>5.87</sub>	23.85 <sub>5.43</sub>	31.49 <sub>1.45</sub>	20.24 <sub>6.56</sub>	19.28 <sub>1.91</sub>	22.77 <sub>3.17</sub>	10.62 <sub>0.84</sub>	43.81 <sub>0.76</sub>	14.43 <sub>2.40</sub>	24.67 <sub>0.92</sub>
GPT-2	$\lambda$	0.014	0.1	0.08	0.06	0.002	0.08	0.2	0.04	0.014	0.04	0.04	0.2	—
	Acc. (%)	Val. w/o	45.44 <sub>1.50</sub>	39.26 <sub>0.95</sub>	44.65 <sub>0.51</sub>	42.83 <sub>0.93</sub>	51.41 <sub>0.76</sub>	69.06 <sub>0.19</sub>	43.20 <sub>0.41</sub>	53.03 <sub>1.35</sub>	59.63 <sub>0.28</sub>	29.67 <sub>0.77</sub>	38.55 <sub>1.43</sub>	41.82 <sub>0.51</sub>
		Val. w/	49.17 <sub>1.43</sub>	64.06 <sub>3.41</sub>	46.05 <sub>0.66</sub>	46.99 <sub>1.77</sub>	52.11 <sub>1.50</sub>	59.65 <sub>3.26</sub>	51.64 <sub>1.67</sub>	58.26 <sub>2.02</sub>	61.50 <sub>1.23</sub>	31.17 <sub>0.94</sub>	42.44 <sub>1.86</sub>	51.48 <sub>4.38</sub>
	Test (%)	w/o	39.62 <sub>1.33</sub>	43.62 <sub>1.17</sub>	40.03 <sub>1.26</sub>	39.41 <sub>1.29</sub>	50.89 <sub>1.04</sub>	63.12 <sub>0.38</sub>	45.93 <sub>0.51</sub>	52.27 <sub>0.81</sub>	59.47 <sub>1.41</sub>	27.43 <sub>0.59</sub>	29.48 <sub>1.29</sub>	41.66 <sub>0.57</sub>
		w/	40.95 <sub>1.42</sub>	67.25 <sub>4.00</sub>	45.28 <sub>3.50</sub>	36.90 <sub>0.18</sub>	51.01 <sub>0.91</sub>	61.20 <sub>0.68</sub>	50.90 <sub>2.63</sub>	56.19 <sub>1.18</sub>	61.54 <sub>1.51</sub>	33.38 <sub>1.16</sub>	34.04 <sub>1.15</sub>	50.28 <sub>2.44</sub>
	MF1 (%)	Val. w/o	25.62 <sub>1.81</sub>	17.88 <sub>0.37</sub>	37.69 <sub>0.08</sub>	37.71 <sub>0.88</sub>	50.82 <sub>0.67</sub>	41.45 <sub>0.49</sub>	34.36 <sub>0.41</sub>	35.02 <sub>0.00</sub>	57.74 <sub>0.42</sub>	20.33 <sub>0.49</sub>	33.74 <sub>1.86</sub>	30.98 <sub>0.50</sub>
		Val. w/	26.05 <sub>1.59</sub>	24.20 <sub>0.85</sub>	32.40 <sub>0.25</sub>	30.78 <sub>1.69</sub>	51.57 <sub>1.59</sub>	48.90 <sub>1.11</sub>	49.38 <sub>0.17</sub>	35.35 <sub>1.00</sub>	61.28 <sub>1.18</sub>	22.14 <sub>0.91</sub>	34.46 <sub>1.13</sub>	46.20 <sub>3.91</sub>
GPT-2	Test (%)	w/o	24.70 <sub>1.17</sub>	18.24 <sub>0.52</sub>	35.78 <sub>1.49</sub>	38.79 <sub>1.25</sub>	49.37 <sub>1.31</sub>	39.90 <sub>0.94</sub>	35.56 <sub>0.75</sub>	35.78 <sub>1.11</sub>	59.46 <sub>1.42</sub>	18.89 <sub>0.62</sub>	29.36 <sub>1.29</sub>	34.43 <sub>0.79</sub>
		w/	24.49 <sub>1.20</sub>	23.74 <sub>1.10</sub>	34.03 <sub>1.20</sub>	27.41 <sub>1.97</sub>	49.56 <sub>0.99</sub>	41.46 <sub>1.26</sub>	47.25 <sub>3.97</sub>	36.67 <sub>2.14</sub>	61.15 <sub>1.50</sub>	24.06 <sub>1.39</sub>	31.48 <sub>3.32</sub>	49.05 <sub>1.35</sub>
	$ECE_1$ (%)	Val. w/o	5.94 <sub>0.81</sub>	36.18 <sub>0.96</sub>	14.79 <sub>0.92</sub>	15.51 <sub>1.25</sub>	30.35 <sub>1.12</sub>	19.39 <sub>0.40</sub>	47.50 <sub>0.37</sub>	8.33 <sub>1.12</sub>	2.10 <sub>0.88</sub>	30.85 <sub>1.06</sub>	14.17 <sub>1.04</sub>	51.18 <sub>0.62</sub>
		Val. w/	5.74 <sub>1.40</sub>	11.64 <sub>2.21</sub>	16.36 <sub>4.41</sub>	17.55 <sub>1.50</sub>	29.51 <sub>1.54</sub>	21.23 <sub>1.51</sub>	24.01 <sub>0.97</sub>	8.20 <sub>1.78</sub>	2.99 <sub>1.10</sub>	30.92 <sub>1.99</sub>	13.52 <sub>2.18</sub>	23.93 <sub>5.48</sub>
	Test (%)	w/o	9.29 <sub>1.37</sub>	28.14 <sub>1.21</sub>	17.90 <sub>1.04</sub>	15.04 <sub>1.29</sub>	31.42 <sub>1.13</sub>	23.02 <sub>0.54</sub>	45.21 <sub>0.78</sub>	8.00 <sub>0.92</sub>	1.89 <sub>0.55</sub>	32.76 <sub>0.80</sub>	22.30 <sub>1.33</sub>	48.69 <sub>0.51</sub>
		w/	8.16 <sub>1.12</sub>	9.01 <sub>2.57</sub>	18.69 <sub>1.81</sub>	25.35 <sub>2.98</sub>	31.07 <sub>1.22</sub>	26.06 <sub>1.64</sub>	25.52 <sub>4.12</sub>	6.88 <sub>1.71</sub>	3.12 <sub>0.80</sub>	29.73 <sub>2.28</sub>	19.15 <sub>4.69</sub>	21.63 <sub>2.11</sub>

## Our contribution can be summarized as:

- We propose NOISYICL, which simply adds noise into LMs and then executes ICL (§2). Our experiment shows that NOISYICL obtains a better ICL performance (§3.2).
- We show that adding noise effectively calibrates LMs to reduce prediction bias and unfaithful confidence in ICL (§3.3).

## 2 NoisyICL

Here we introduce the basic form of ICL and our perturbation method named NOISYICL.

**In-context Learning.** Suppose a supervised classification dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ , where  $x_i$  is an input, and  $y_i \in \mathbb{U}$  is the corresponding label in a label space  $\mathbb{U}$ . For each query  $x_q$  to be predicted, we sample a demo sequence  $G = \{(x_{a_j}, y_{a_j})\}_{j=1}^k$  from  $\mathcal{D}$ , where  $k$  is the number of demos, and construct a prompt input in natural language form  $s = f(G, x_q)$  with a pattern  $f$ . Then, we input  $s$  into the LM  $P_\theta(\cdot)$  with parameters  $\theta$  and get an output token distribution  $P_\theta(\cdot|s)$ . We choose the label token  $l$  with the maximum probability **among the label space** as the prediction  $\hat{y}_q$ , that is:

$$\hat{y}_q = \operatorname{argmax}_{l \in \mathbb{U}} P_\theta(l|s). \quad (1)$$

Notice that we only construct prompts to drive the model to predict answers generatively, without any parameter updates.

**NOISYICL.** For each parameter matrix  $\theta_i$  in the LM  $P_\theta(\cdot)$  used for ICL, we simply do an interpolation between  $\theta_i$  and a noise matrix sampled from an isotropic normal distribution  $\mathcal{N}(0, \sigma^2)$  with intensity  $\lambda$ , that is:

$$\theta'_i = (1 - \lambda)\theta_i + \lambda\mathcal{N}(0, \sigma^2), \quad (2)$$

where  $\lambda$  and  $\sigma$  are model or task-specific hyperparameters. Then we perform the aforementioned ICL with the interpolated LM  $P_{\theta'}(\cdot)$ .

## 3 Experiments and Results

We investigate the effectiveness and calibration abilities of NOISYICL. We find that NOISYICL can improve ICL performance by an average of 10% (§3.2) and effectively calibrate the model's prediction bias and unfaithful confidence (§3.3).

### 3.1 Settings

**Data.** We use 12 commonly used classification datasets in our experiments. Since some of the datasets do not provide valid splitting, we randomly split these datasets into validation sets and test sets (see Appendix A for further details).

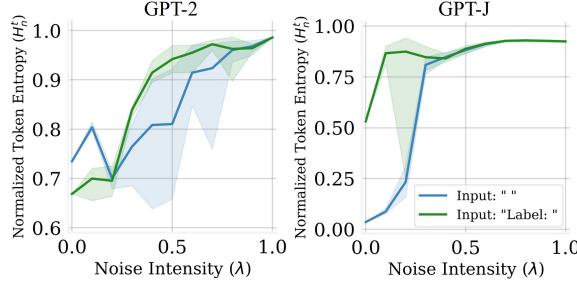


Figure 2: The correlation between the normalized token entropy  $H_n^t$  and the noise intensity  $\lambda$  with empty inputs. When the noise gets stronger, the  $H_n^t$  becomes higher, which indicates a fairer output.

**Models.** We use GPT-2 (parameters: 137M) (Radford et al., 2019) and GPT-J (parameters: 6053M) (Wang and Komatsuzaki, 2021). The model checkpoints are loaded from huggingface<sup>2</sup>.

**Hyperparameters.** We fix  $\sigma$ , the standard deviation of the noise distribution, to 0.02, which is the same as the initialization distribution of both models. We believe the noise intensity  $\lambda$  should vary w.r.t. datasets and models. Therefore, we determine the most suitable  $\lambda$  by a simple search method for each dataset and model on the validation set (details in Appendix D). The selected intensities are shown in Table 1, which are concentrated in  $(0, 0.3]$ . Moreover, we confirm that the  $\lambda$  keeps relatively stable w.r.t.  $k$  given a dataset and model (Appendix D).

**Other details.** We default to use four demos and a simple prompt template as shown in Appendix B. Every labeled data is treated as the test query once, and for each query, we do 2 tries. Each experiment is repeated 10 times with different noise matrices.

### 3.2 NOISYICL Improves ICL Performance

We show accuracy and macro-F1 on the 12 downstream datasets with and without NOISYICL on the optimal  $\lambda$  in Table 1, and averaged results in Fig. 1 (see Appendix C for results of other numbers of demos ( $k$ )).

The results show that NOISYICL produces an average performance improvement of around 10%. This phenomenon preliminarily confirms our hypothesis: NOISYICL fits the pre-trained LMs towards ICL.

However, such gains vary depending on the dataset and model. In some combinations of

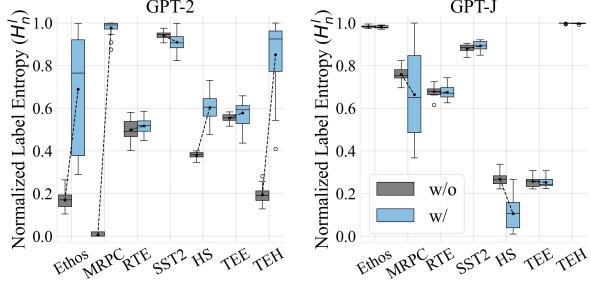


Figure 3: The normalized label entropy  $H_n^l$  on both models and 7 datasets with and without appropriate-noised NOISYICL. In most cases, the  $H_n^l$  with NOISYICL (w/) is greater than without NOISYICL (w/o).

datasets and models, significant performance improvement cannot be observed. We speculate the reason is that NOISYICL does not provide new knowledge from new training examples, and these datasets are too hard for the model intrinsically no matter whether NOISYICL is used.

### 3.3 NOISYICL Is A Calibration

This section finds that NOISYICL conducts the following aspects of calibrations:

- **Lower prediction bias.** When no valid query is given, the predicted labels should have balanced frequencies. However, predictions conducted by under-calibrated LMs usually have significant bias, which is harmful to ICL (Fei et al., 2023; Zhao et al., 2021; Wei et al., 2023; Lu et al., 2022). Eliminating these biases can be seen as an aspect of calibration.
- **More faithful confidence.** In classification tasks, the outputted probability of the predicted label is called confidence. Suitable confidence should faithfully reflect the accuracy of outputs, that is, a prediction with greater accuracy should be assigned with greater confidence (Corbière et al., 2019). It has been proven that faithful confidence improves the stability of the model (Guo et al., 2017; Grabinski et al., 2022). Making the confidence more faithful is also an aspect of calibration (Guo et al., 2017; Tian et al., 2023).

Specifically, we find that NOISYICL induces the model to predict labels with more fair and faithful confidence.

**NOISYICL reduces prediction bias.** We investigate the following two types of prediction bias. (See Appendix E for further calculation details.)

<sup>2</sup>[huggingface.co/gpt2](https://huggingface.co/gpt2), and [huggingface.co/EleutherAI/gpt-j-6b](https://huggingface.co/EleutherAI/gpt-j-6b)

177           **(1) Intrinsic token bias.** LMs are pre-trained  
 178 with natural corpus, where different tokens have  
 179 various frequencies. LMs learn such frequencies  
 180 and act as biases among different tokens in pre-  
 181 diction, which are deemed harmful to ICL (Fei  
 182 et al., 2023). We believe that such token bias can  
 183 be reduced by NOISYICL.

184 To quantify this, we calculate normalized token  
 185 entropy  $H_n^t$ : the entropy of the predicted probabili-  
 186 ty distribution among the whole vocabulary given  
 187 an empty input (we use “”, and “Label: ”), w.r.t.  
 188 various noise intensities.

189 The results are shown in Fig. 2. A clear positive  
 190 correlation between the noise intensity and  $H_n^t$   
 191 can be observed, meaning the model gives a fairer  
 192 output when more noise is given.

193           **(2) Overall label bias.** In classification tasks,  
 194 when no query is given, equal probabilities should  
 195 be assigned to labels fairly, which benefits ICL (Lu  
 196 et al., 2022). We believe that NOISYICL promotes  
 197 such fairness.

198 To quantify this, we use normalized label entro-  
 199 py  $H_n^l$ : the entropy of the predicted probability  
 200 distribution among the label space given some de-  
 201 mos and an empty query. Notice that it is natural  
 202 for a model to predict a “neutral” label when no  
 203 query is given, so label bias on a dataset with a  
 204 “neutral” label cannot be fairly measured with  $H_n^l$ .  
 205 Therefore, we ignore these datasets with “neutral”  
 206 labels (details in Appendix A).

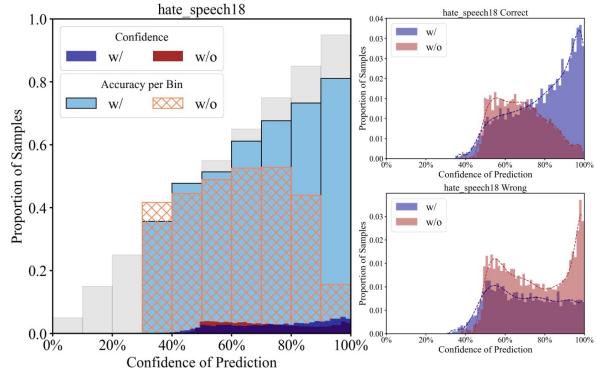
207 The results are shown in Fig. 3, indicating that in  
 208 most cases, NOISYICL calibrates the overall bias,  
 209 especially on GPT-2.

210           **NOISYICL promotes faithful confidence.** The  
 211 Expected Calibration Error ( $ECE_p$ ) (Naeini et al.,  
 212 2015) is a widely-used indicator for the faithfulness  
 213 of confidence in classification tasks:

$$ECE_p = \mathbb{E}(|\max(\hat{z}) - \mathbb{E}_{y=\text{argmax}_i \hat{z}_i}(1)|^p)^{\frac{1}{p}}, \quad (3)$$

215 where  $\hat{z}$  is the predicted probability vector on the  
 216 label space by a classification model, the final pre-  
 217 diction ( $\text{argmax}_i \hat{z}_i$ ) can be obtained with a confi-  
 218 dence ( $\max \hat{z}$ ), and the ground-truth label is  $y$ .

219 Let  $p = 1$ , we use the  $ECE_1$  to investigate  
 220 the faithfulness of the ICL output. The details  
 221 of the calculation are shown in Appendix F. A  
 222 lower  $ECE_1$  means more faithful confidence, that  
 223 is, the confidence becomes a more accurate pre-  
 224 diction of accuracy (Corbière et al., 2019). Fig. 4  
 225 shows a case study of GPT-2 on the hate\_speech18



226           **Figure 4: Left:** Reliability diagrams (sparse bars) and  
 227 global confidence distribution (dense bars) of GPT-2 on  
 228 hate\_speech18 with (**w/**,  $ECE_1 = 9.01\%$ ) and without  
 229 (**w/o**,  $ECE_1 = 28.14\%$ ) NOISYICL. The predictions  
 230 are divided into bins according to confidence, and we  
 231 visualize the accuracy of each bin as a histogram. The grey  
 232 bars are ideal, which is closer to the one with NOISY-  
 233 ICL. **Right upper:** Confidence distribution on correct  
 234 predictions. Relatively right-shifted with NOISYICL.  
 235 **Right lower:** Confidence distribution on wrong predictions.  
 236 Relatively left-shifted with NOISYICL. ( $k = 4$ )

237 dataset, indicating that the prediction is more faithful-  
 238 ful with NOISYICL (full visualized data are in  
 239 Appendix G). We show the results with and without  
 240 appropriate-noised NOISYICL for  $ECE_1$  on  
 241 the 12 datasets in Table 1.

242 In most cases, the  $ECE_1$  is reduced by around  
 243 25% by NOISYICL, which means the confidence is  
 244 more faithful with NOISYICL. This suggests that  
 245 NOISYICL can drive the model to produce more  
 246 faithful confidence, that is, less over-confidence  
 247 in wrong predictions, and less under-confidence  
 248 in correct predictions, which suggests that NOISYICL  
 249 solves the confidence calibration.

250 Based on the above two investigations, we be-  
 251 lieve that NOISYICL is an effective calibration for  
 252 the ICL scenario. This can be a reasonable expla-  
 253 nation for the observed performance improvement  
 254 in §3.2.

## 4 Conclusion

255 In this paper, we propose NOISYICL, which simply  
 256 adds random noise to the parameters of LMs to fit  
 257 them from pre-trained knowledge to ICL. We show  
 258 that NOISYICL can improve the ICL performance  
 259 and also calibrate the model for fairer outputs and  
 260 more faithful confidence.

## 251 5 Limitations

252 As mentioned before, unlike the fine-tuning on ad-  
253 ditional ICL-style datasets (Min et al., 2022; Zhao  
254 et al., 2021; Wei et al., 2021), NOISYICL does  
255 not provide new knowledge for the model, so the  
256 calibrated model can not discover tasks that are not  
257 potentially included in the pre-training data (Gu  
258 et al., 2023). Moreover, a naive search for the best  
259 noise intensity is not efficient and satisfactory. An  
260 effective detection of  $\lambda$  should be developed. Cur-  
261 rently, we can only confirm that for one dataset  
262 and model, the  $\lambda$  keeps relatively stable, especially  
263 when a large  $k$  is given (Appendix. D).

264 **Future Works.** Besides fixing the limits, future  
265 works can focus on where and how the noise should  
266 be introduced. In Transformer-based models, dif-  
267 ferent layers have different abilities (Wang et al.,  
268 2023; Jawahar et al., 2019; Kobayashi et al., 2023).  
269 So, treating these layers differently may be an ef-  
270 fective improvement of NOISYICL. An isotropic  
271 normal distribution is too simple, noise sampling  
272 methods should be discussed.

273 Moreover, adding noise to model parameters can  
274 be an erasing of pre-training (Ilharco et al., 2022),  
275 so, the search for  $\lambda$  is the search for the best check-  
276 point of pre-training. With these checkpoints, we  
277 can determine (Han et al., 2023; Han and Tsvetkov,  
278 2022) which data is disadvantageous to ICL, and  
279 what knowledge is essential to ICL, to better reveal  
280 the essence of ICL.

## 281 References

282 Valerio Basile, Cristina Bosco, Elisabetta Fersini,  
283 Debora Nozza, Viviana Patti, Francisco Manuel  
284 Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti.  
285 2019. SemEval-2019 task 5: Multilingual detection  
286 of hate speech against immigrants and women in  
287 Twitter. In *Proceedings of the 13th International*  
288 *Workshop on Semantic Evaluation*, pages 54–63, Min-  
289 neapolis, Minnesota, USA. Association for Compu-  
290 tational Linguistics.

291 Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo  
292 Giampiccolo. 2009. The fifth pascal recognising  
293 textual entailment challenge. *TAC*, 7:8.

294 Charles Corbière, Nicolas Thome, Avner Bar-Hen,  
295 Matthieu Cord, and Patrick Pérez. 2019. Address-  
296 ing failure prediction by learning model confidence.  
297 *Advances in Neural Information Processing Systems*,  
298 32.

299 Ido Dagan, Oren Glickman, and Bernardo Magnini.  
300 2005. The pascal recognising textual entailment chal-

lenge. In *Machine learning challenges workshop*,  
301 pages 177–190. Springer.

302 Ona de Gibert, Naiara Perez, Aitor García-Pablos, and  
303 Montse Cuadros. 2018. Hate Speech Dataset from  
304 a White Supremacy Forum. In *Proceedings of the*  
305 *2nd Workshop on Abusive Language Online (ALW2)*,  
306 pages 11–20, Brussels, Belgium. Association for  
307 Computational Linguistics.

308 Bill Dolan and Chris Brockett. 2005. Automatically  
309 constructing a corpus of sentential paraphrases.  
310 In *Third International Workshop on Paraphrasing*  
311 (*IWP2005*).

312 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiy-  
313 ong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and  
314 Zhifang Sui. 2022. A survey for in-context learning.  
315 *arXiv preprint arXiv:2301.00234*.

316 Yu Fei, Yifan Hou, Zeming Chen, and Antoine Bosselut.  
317 2023. Mitigating label biases for in-context learning.  
318 In *Proceedings of the 61st Annual Meeting of the*  
319 *Association for Computational Linguistics (Volume 1:*  
320 *Long Papers*), pages 14014–14031, Toronto, Canada.  
321 Association for Computational Linguistics.

322 Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and  
323 William B Dolan. 2007. The third pascal recognizing  
324 textual entailment challenge. In *Proceedings of the*  
325 *ACL-PASCAL workshop on textual entailment and*  
326 *paraphrasing*, pages 1–9.

327 Julia Grabinski, Paul Gavrikov, Janis Keuper, and Mar-  
328 gret Keuper. 2022. Robust models are less over-  
329 confident. *Advances in Neural Information Process-  
330 ing Systems*, 35:39059–39075.

331 Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang.  
332 2023. Pre-training to learn in context. *arXiv preprint*  
333 *arXiv:2305.09137*.

334 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Wein-  
335 berger. 2017. On calibration of modern neural net-  
336 works. In *International conference on machine learn-  
337 ing*, pages 1321–1330. PMLR.

338 R Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo  
339 Giampiccolo, Bernardo Magnini, and Idan Szpektor.  
340 2006. The second pascal recognising textual entail-  
341 ment challenge. In *Proceedings of the Second PAS-  
342 CAL Challenges Workshop on Recognising Textual  
343 Entailment*, volume 7, pages 785–794.

344 Xiaochuang Han, Daniel Simig, Todor Mihaylov, Yu-  
345 lia Tsvetkov, Asli Celikyilmaz, and Tianlu Wang.  
346 2023. Understanding in-context learning via sup-  
347 portive pretraining data. In *Proceedings of the 61st*  
348 *Annual Meeting of the Association for Computational*  
349 *Linguistics (Volume 1: Long Papers)*, pages 12660–  
350 12673.

351 Xiaochuang Han and Yulia Tsvetkov. 2022. Orca: In-  
352 terpreting prompted language models via locating  
353 supporting data evidence in the ocean of pretraining  
354 data. *arXiv preprint arXiv:2205.12600*.

356	Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2022. Editing models with task arithmetic.	Maria Pontiki, Dimitris Galanis, John Pavlopoulos, Harris Papageorgiou, Ion Androutsopoulos, and Suresh Manandhar. 2014. <i>SemEval-2014 task 4: Aspect based sentiment analysis</i> . In <i>Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)</i> , pages 27–35, Dublin, Ireland. Association for Computational Linguistics.	409 410 411 412 413 414 415
357	In <i>The Eleventh International Conference on Learning Representations</i> .		
358			
359			
360			
361	Ganesh Jawahar, Benoît Sagot, and Djamé Seddah.	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	416 417 418
362	2019. What does bert learn about the structure of language? In <i>ACL 2019-57th Annual Meeting of the Association for Computational Linguistics</i> .		
363			
364			
365	Goro Kobayashi, Tatsuki Kurabayashi, Sho Yokoi, and Kentaro Inui. 2023. <i>Transformer language models handle word frequency in prediction head</i> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 4523–4535, Toronto, Canada. Association for Computational Linguistics.	Sara Rosenthal, Noura Farra, and Preslav Nakov. 2017. <i>Semeval-2017 task 4: Sentiment analysis in twitter</i> . In <i>Proceedings of the 11th international workshop on semantic evaluation (SemEval-2017)</i> , pages 502–518.	419 420 421 422 423
366			
367			
368			
369			
370			
371	Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8086–8098.	Emily Sheng and David Uthus. 2020. Investigating societal biases in a poetry composition system. <i>arXiv preprint arXiv:2011.02686</i> .	424 425 426
372			
373			
374			
375			
376			
377			
378	P. Malo, A. Sinha, P. Korhonen, J. Wallenius, and P. Takala. 2014. Good debt or bad debt: Detecting semantic orientations in economic texts. <i>Journal of the Association for Information Science and Technology</i> , 65.	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.	427 428 429 430 431 432 433
379			
380			
381			
382			
383	Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2022. Metacl: Learning to learn in context. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2791–2809.	Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. <i>arXiv preprint arXiv:2305.14975</i> .	434 435 436 437 438 439
384			
385			
386			
387			
388			
389	Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. Semeval-2018 task 1: Affect in tweets. In <i>Proceedings of the 12th international workshop on semantic evaluation</i> , pages 1–17.	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. In <i>International Conference on Learning Representations</i> .	440 441 442 443 444
390			
391			
392			
393			
394	Ioannis Mollas, Zoe Chrysopoulou, Stamatis Karlos, and Grigoris Tsoumakas. 2020. <i>Ethos: an online hate speech detection dataset</i> .	Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <a href="https://github.com/kingoflolz/mesh-transformer-jax">https://github.com/kingoflolz/mesh-transformer-jax</a> .	445 446 447 448
395			
396			
397	Marius Mosbach, Tiago Pimentel, Shauli Ravfogel, Dietrich Klakow, and Yanai Elazar. 2023. Few-shot fine-tuning vs. in-context learning: A fair comparison and evaluation. <i>arXiv preprint arXiv:2305.16938</i> .	Lean Wang, Lei Li, Damai Dai, Deli Chen, Hao Zhou, Fandong Meng, Jie Zhou, and Xu Sun. 2023. Label words are anchors: An information flow perspective for understanding in-context learning. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9840–9855, Singapore. Association for Computational Linguistics.	449 450 451 452 453 454 455 456
398			
399			
400			
401	Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. 2015. Obtaining well calibrated probabilities using bayesian binning. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 29.	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. In <i>International Conference on Learning Representations</i> .	457 458 459 460 461
402			
403			
404			
405	Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. 2023. What in-context learning "learns" in-context: Disentangling task recognition and task learning. <i>arXiv preprint arXiv:2305.09731</i> .	Jerry Wei, Le Hou, Andrew Lampinen, Xiangning Chen, Da Huang, Yi Tay, Xinyun Chen, Yifeng Lu, Denny	462 463
406			
407			
408			

464 Zhou, Tengyu Ma, et al. 2023. Symbol tuning im-  
 465 proves in-context learning in language models. *arXiv*  
 466 preprint arXiv:2305.08298.

467 Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng  
 468 Huang. 2022. Noisytune: A little noise can help  
 469 you finetune pretrained language models better. In  
 470 *Proceedings of the 60th Annual Meeting of the As-*  
 471 *sociation for Computational Linguistics (Volume 2:*  
 472 *Short Papers)*, pages 680–685.

473 Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and  
 474 Sameer Singh. 2021. Calibrate before use: Improv-  
 475 ing few-shot performance of language models. In *In-*  
 476 *International Conference on Machine Learning*, pages  
 477 12697–12706. PMLR.

## A Datasets

479 The datasets used in this paper are shown in Ta-  
 480 ble 2.

481 \*To construct inputs of appropriate length, we re-  
 482 move data with a length exceeding 500 from the  
 483 Ethos, and the number of the remaining data is 980.

484 **Dataset Splitting.** For each dataset, we first shuf-  
 485 fle it with random seed 42. Then, we choose the  
 486 512 data at the tail as the testing data, and the 512  
 487 data at the head (if the total number of data is less  
 488 than 1024, we choose the rest of the data.) as the  
 489 validation data.

## B Prompt Patterns

491 In this paper, we use a minimum prompt template.  
 492 For each task category, we design a template as  
 493 shown below.

494 For single-sentence classification datasets  $(x, y)$ ,  
 495 we use:

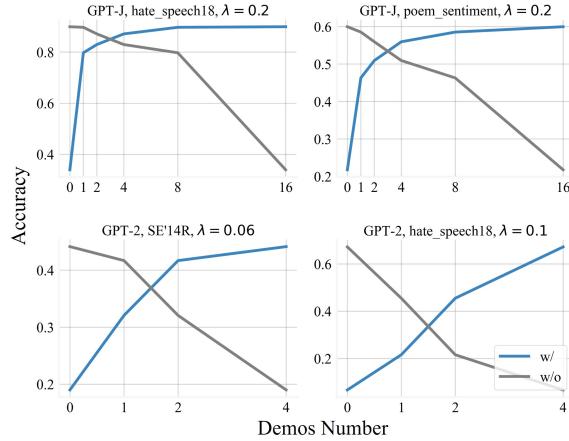
496 Input: <x>, Label: <y> \n  
 497 ...  
 498 Input: <x>, Label:

499 For aspect-based sentiment classification  
 500 datasets  $((x, a), y)$ , we use:

501 Input: <x>, Aspect: <a>, Label: <y> \n  
 502 ...  
 503 Input: <x>, Aspect: <a>, Label:

504 For double-sentence classification datasets  
 505  $((x_1, x_2), y)$ . we use:

506 Input: <x1>, Text 2: <x2>, Label: <y> \n  
 507 ...  
 508 Input: <x1>, Text 2: <x2>, Label:



473 Figure 5: The relationship between demos quantity and  
 474 accuracy in some cases. NOISYICL can make the model  
 475 learn from the demos correctly.

Input: This is delicious!	Label: Positive
Input: I don't like it.	Label: Negative
Input: Boring.	Label: Negative
Input: It is fantastic!	Label: Positive
Input:	Label:

473 Figure 6: An example of inputs in the normalized label  
 474 entropy calculation.

## C Complete Results: Accuracy, 509 Macro-F1, and ECE<sub>1</sub> With Various 510 Demo Numbers ( $k$ )

511 We use various numbers of demos in the experi-  
 512 ments shown in the §3.2.

513 The results of zero-shot ( $k = 0$ ) are shown in Ta-  
 514 ble 3.

515 The results of 1-shot ( $k = 1$ ) are shown in Table 4.  
 516 The results of 2-shot ( $k = 2$ ) are shown in Table 5.  
 517 The results of 8-shot ( $k = 8$ ) are shown in Table 6.  
 518 The results of 16-shot ( $k = 16$ ) are shown in Ta-  
 519 ble 7.

520 Notice that in Table 6 and Table 7, some ex-  
 521 periments are not feasible due to the length of se-  
 522 quences.

523 In these results, NOISYICL still outperforms the  
 524 baseline, which proves that NOISYICL is generally  
 525 applicable with various  $k$ .

## D Searching and Stability of $\lambda$

527 We select a set of candidates of  $\lambda$  as  $\{0, 0.002,$   
 528  $0.004, 0.006, 0.008, 0.01, 0.012, 0.014, 0.016,$   
 529  $0.018, 0.02, 0.04, 0.06, 0.08, 0.1, 0.2, 0.3, 0.4,$   
 530  $0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$ , test the performance of

Table 2: Datasets and Abbreviations used in this paper. **Abbr.:** Abbreviation, **Data#:** The total number of data, **Label#:** The number of classes, **Neutral?:** Does the dataset have a neutral label?, **Major.**  $\sim$  **Minor.(%)**: The proportion of majority labels  $\sim$  the proportion of minority labels.

Dataset	Abbr.	Data#	Label#	Neutral?	Major.	$\sim$ Minor.(%)
<i>single-sentence classification:</i>						
poem_sentiment (Sheng and Uthus, 2020)	PS	1101	4	✓	62.2	$\sim$ 5.5
hate_speech18 (de Gibert et al., 2018)	HS	10944	4	✗	86.9	$\sim$ 1.5
Ethos* (binary) (Mollas et al., 2020)	—	980	2	✗	56.6	$\sim$ 43.4
financial_phrasebank (all agree) (Malo et al., 2014)	FP	2264	3	✓	61.4	$\sim$ 13.4
GLUE-SST2 (Wang et al., 2018; Socher et al., 2013)	SST2	68221	2	✗	55.8	$\sim$ 44.2
tweet_eval_emotion (Mohammad et al., 2018)	TEE	5052	4	✗	43	$\sim$ 9
tweet_eval_sentiment (Rosenthal et al., 2017)	TES	59899	3	✓	45.3	$\sim$ 15.5
tweet_eval_hate (Basile et al., 2019)	TEH	12970	2	✗	58	$\sim$ 42
<i>aspect-based sentiment classification:</i>						
SemEval 2014-Task 4 Restaurants (Pontiki et al., 2014)	SE'14R	4722	3	✓	61.2	$\sim$ 17.6
SemEval 2014-Task 4 Laptops (Pontiki et al., 2014)	SE'14L	2951	3	✓	45.3	$\sim$ 15.5
<i>double-sentence classification:</i>						
GLUE-RTE (Wang et al., 2018; Dagan et al., 2005) (Haim et al., 2006; Giampiccolo et al., 2007; Bentivogli et al., 2009)	RTE	2767	2	✗	50.2	$\sim$ 49.8
GLUE-MRPC (Wang et al., 2018; Dolan and Brockett, 2005)	MRPC	4076	2	✗	67.4	$\sim$ 32.6

NOISYICL with this set of  $\lambda$ , and select the one with the best accuracy as the optimal  $\lambda$ .

Moreover, we find that given a dataset and model,  $\lambda$  remains relatively stable w.r.t the demo number  $k$ , especially when  $k$  is large. As shown in Fig. 7 and Fig. 8, we calculate  $|\lambda_i - \lambda_j|$ , the distance between the optimal  $\lambda$  of different demo numbers as the heatmap, and the normalized remaining range  $R(\lambda_{i:}) / \max(\lambda)$ , where the  $R(\lambda_{i:})$  is the range of all the optimal  $\lambda_c$  where the demo number  $c \geq i$ , and the  $\max(\lambda)$  is the maximum of all the optimal  $\lambda$  on the dataset. It characterizes the dispersion of  $\lambda$  when the demo number is greater than the given one  $i$ .

The results show that while the  $k$  is increasing, the distance and the normalized remaining range quickly zero out, which supports our hypothesis: NOISYICL is a bridge from the pre-training to ICL. For that during this process, the demo number should have a minimal impact.

## E Calculation of $H_n^t$ and $H_n^l$

**Normalized Token Entropy.** We calculate the normalized token entropy  $H_n^t$  of the predicted probability distribution of tokens with an empty input as:

$$H_n^t = -\frac{\sum_{l \in \mathbb{V}} P_\theta(l|x_\emptyset) \ln P_\theta(l|x_\emptyset)}{\ln |\mathbb{V}|}, \quad (4)$$

where the  $P_\theta(\cdot)$  is an LM with a vocabulary space  $\mathbb{V}$  and size  $|\mathbb{V}|$ , the  $x_\emptyset$  is an empty input, such as “ ”, or “Label: ”. The model outputs a global probability distribution of tokens when the  $x_\emptyset$  is

given, and we calculate the normalized entropy on the distribution.

**Normalized Label Entropy.** We calculate the normalized label entropy  $H_n^l$  of the model prediction probability distribution on the label space with 4 demos and an empty queue from a dataset (Fig. 6 is an example) as:

$$H_n^l = -\frac{\sum_{l \in \mathbb{U}} r_{P,\theta}(l|x_*) \ln r_{P,\theta}(l|x_*)}{\ln |\mathbb{U}|}, \quad (5)$$

where the  $r_{P,\theta}(l|x_*)$  is the predicted frequency of label  $l$  given the aforementioned query-less only-demos input  $x_*$ . We use 512 tries to estimate this frequency. The  $\mathbb{U}$  and  $|\mathbb{U}|$  are the label space and label space size, respectively.

## F Calculation of $ECE_1$

Given a prediction set  $\mathcal{Y} = \{(\hat{y}_i, z_i, y_i)\}_{i=1}^n$  of predictions  $\hat{y}_i$ , predicted confidence  $z_i \in [0, 1]$ , and supervised label  $y_i$ , in the calculation of  $ECE_1$ , we first divide the confidence space  $(0, 1)$  into  $m = 10$  equal bins  $B_j = [0.1(j-1), 0.1j]_{j=1}^m$  according to the  $z_i$  as shown in Fig. 4. The amount of prediction in each bin is denoted as  $|B_j|$ . For each bin, we calculate the accuracy  $\alpha_j = \frac{1}{|B_j|} \sum_{\hat{y}_i \in B_j, \hat{y}_i = y_1} 1$  of predictions in the  $j$ th bin, and assign a standard accuracy based on the average of confidence in the bin, that is,  $\alpha_j^s = \frac{1}{|B_j|} \sum_{z_i \in B_j} z_i$ . The  $ECE_1$  can be described as:

$$ECE_1 = \sum_{j=1}^m \frac{|B_j|}{n} |\alpha_j - \alpha_j^s|, \quad (6)$$

Table 3: Accuracy and Macro-F1 results (% ,  $mean_{std}$ ,  $k = 0$ ). A better result is in green. Notation is the same with the Tabel 1.

Dataset		PS	HS	SE'14R	SE'14L	RTE	MRPC	Ethos	FP	SST2	TEE	TES	TEH	Mean
$\lambda$		0.1	0.1	0.1	0.1	0.06	0.02	0.004	0.02	0.1	0.2	0.008	0.3	—
Acc. (%)	Val.	w/o	16.60 <sub>0.00</sub>	37.30 <sub>0.00</sub>	33.98 <sub>0.00</sub>	38.87 <sub>0.00</sub>	52.54 <sub>0.00</sub>	65.82 <sub>0.00</sub>	51.07 <sub>0.00</sub>	61.13 <sub>0.00</sub>	57.42 <sub>0.00</sub>	39.84 <sub>0.00</sub>	49.22 <sub>0.00</sub>	43.36 <sub>0.00</sub>
		w/	21.86 <sub>6.75</sub>	50.70 <sub>4.33</sub>	40.74 <sub>10.27</sub>	45.94 <sub>2.72</sub>	53.40 <sub>0.55</sub>	67.77 <sub>0.55</sub>	51.07 <sub>0.19</sub>	61.37 <sub>0.60</sub>	60.98 <sub>3.87</sub>	39.57 <sub>3.50</sub>	49.53 <sub>0.20</sub>	54.14 <sub>5.91</sub>
	Test	w/o	17.38 <sub>0.00</sub>	29.88 <sub>0.00</sub>	28.91 <sub>0.00</sub>	39.45 <sub>0.00</sub>	49.80 <sub>0.00</sub>	59.18 <sub>0.00</sub>	54.10 <sub>0.00</sub>	61.91 <sub>0.00</sub>	66.41 <sub>0.00</sub>	41.60 <sub>0.00</sub>	49.22 <sub>0.00</sub>	46.88 <sub>0.00</sub>
GPT-J	MF1 (%)	Val.	w/o	11.32 <sub>0.00</sub>	20.81 <sub>0.00</sub>	35.19 <sub>0.00</sub>	36.56 <sub>0.00</sub>	46.42 <sub>0.00</sub>	48.13 <sub>0.00</sub>	48.58 <sub>0.00</sub>	36.13 <sub>0.00</sub>	56.98 <sub>0.00</sub>	15.99 <sub>0.00</sub>	29.81 <sub>0.00</sub>
		w/	14.87 <sub>2.10</sub>	25.82 <sub>2.65</sub>	39.90 <sub>8.62</sub>	41.36 <sub>3.98</sub>	50.10 <sub>1.31</sub>	49.08 <sub>0.54</sub>	48.58 <sub>0.25</sub>	37.35 <sub>0.74</sub>	59.80 <sub>4.50</sub>	25.30 <sub>2.60</sub>	30.36 <sub>0.60</sub>	39.66 <sub>6.23</sub>
	Test	w/o	12.70 <sub>0.00</sub>	16.02 <sub>0.00</sub>	27.01 <sub>0.00</sub>	39.38 <sub>0.00</sub>	45.56 <sub>0.00</sub>	40.48 <sub>0.00</sub>	50.99 <sub>0.00</sub>	34.91 <sub>0.00</sub>	64.45 <sub>0.00</sub>	16.26 <sub>0.00</sub>	29.81 <sub>0.00</sub>	42.58 <sub>0.00</sub>
ECE <sub>1</sub> (% , ↓)	Val.	w/o	36.56 <sub>0.00</sub>	22.03 <sub>0.00</sub>	28.11 <sub>0.00</sub>	17.80 <sub>0.00</sub>	29.83 <sub>0.00</sub>	7.24 <sub>0.00</sub>	28.08 <sub>0.00</sub>	13.45 <sub>0.00</sub>	5.45 <sub>0.00</sub>	48.09 <sub>0.00</sub>	22.10 <sub>0.00</sub>	46.21 <sub>0.00</sub>
		w/	22.52 <sub>6.40</sub>	16.20 <sub>3.75</sub>	16.77 <sub>9.48</sub>	8.39 <sub>3.82</sub>	23.90 <sub>1.51</sub>	9.20 <sub>0.54</sub>	28.12 <sub>0.24</sub>	12.75 <sub>0.37</sub>	4.22 <sub>2.43</sub>	18.66 <sub>2.88</sub>	22.53 <sub>0.19</sub>	21.99 <sub>12.15</sub>
	Test	w/o	35.33 <sub>0.00</sub>	26.92 <sub>0.00</sub>	33.41 <sub>0.00</sub>	16.92 <sub>0.00</sub>	31.46 <sub>0.00</sub>	13.80 <sub>0.00</sub>	24.56 <sub>0.00</sub>	15.09 <sub>0.00</sub>	3.49 <sub>0.00</sub>	46.49 <sub>0.00</sub>	22.10 <sub>0.00</sub>	36.32 <sub>0.00</sub>
GPT-2	MF1 (%)	Val.	w/o	4.30 <sub>0.00</sub>	15.43 <sub>0.00</sub>	23.44 <sub>0.00</sub>	39.84 <sub>0.00</sub>	49.41 <sub>0.00</sub>	69.53 <sub>0.00</sub>	44.23 <sub>0.00</sub>	43.55 <sub>0.00</sub>	60.55 <sub>0.00</sub>	32.03 <sub>0.00</sub>	46.68 <sub>0.00</sub>
		w/	5.00 <sub>0.50</sub>	15.23 <sub>0.00</sub>	40.90 <sub>15.99</sub>	40.35 <sub>5.80</sub>	50.82 <sub>0.42</sub>	63.52 <sub>2.09</sub>	55.98 <sub>1.93</sub>	58.13 <sub>2.45</sub>	63.87 <sub>0.96</sub>	38.71 <sub>0.57</sub>	46.41 <sub>0.16</sub>	49.38 <sub>3.75</sub>
	Test	w/o	5.86 <sub>0.00</sub>	10.74 <sub>0.00</sub>	22.27 <sub>0.00</sub>	35.35 <sub>0.00</sub>	48.63 <sub>0.00</sub>	63.48 <sub>0.00</sub>	46.88 <sub>0.00</sub>	45.51 <sub>0.00</sub>	64.84 <sub>0.00</sub>	31.25 <sub>0.00</sub>	37.50 <sub>0.00</sub>	40.23 <sub>0.00</sub>
ECE <sub>1</sub> (% , ↓)	Val.	w/o	4.70 <sub>0.00</sub>	9.58 <sub>0.00</sub>	24.11 <sub>0.00</sub>	36.57 <sub>0.00</sub>	47.91 <sub>0.00</sub>	41.63 <sub>0.00</sub>	34.16 <sub>0.00</sub>	26.03 <sub>0.00</sub>	54.25 <sub>0.00</sub>	22.91 <sub>0.00</sub>	46.68 <sub>0.00</sub>	29.85 <sub>0.00</sub>
		w/	5.25 <sub>0.57</sub>	9.46 <sub>0.01</sub>	27.57 <sub>5.24</sub>	31.40 <sub>3.63</sub>	50.31 <sub>0.61</sub>	50.06 <sub>1.78</sub>	51.44 <sub>2.19</sub>	33.29 <sub>0.25</sub>	61.13 <sub>2.01</sub>	27.20 <sub>0.36</sub>	38.58 <sub>0.10</sub>	41.68 <sub>2.24</sub>
	Test	w/o	6.08 <sub>0.00</sub>	6.50 <sub>0.00</sub>	20.45 <sub>0.00</sub>	34.84 <sub>0.00</sub>	46.00 <sub>0.00</sub>	38.83 <sub>0.00</sub>	35.37 <sub>0.00</sub>	28.04 <sub>0.00</sub>	64.67 <sub>0.00</sub>	22.73 <sub>0.00</sub>	32.56 <sub>0.00</sub>	29.22 <sub>0.00</sub>
GPT-2	MF1 (%)	Val.	w/o	53.30 <sub>0.00</sub>	46.81 <sub>0.00</sub>	33.37 <sub>0.00</sub>	13.58 <sub>0.00</sub>	18.57 <sub>0.00</sub>	22.00 <sub>0.00</sub>	42.42 <sub>0.00</sub>	13.50 <sub>0.00</sub>	6.37 <sub>0.00</sub>	20.26 <sub>0.00</sub>	5.92 <sub>0.00</sub>
		w/	58.20 <sub>2.84</sub>	48.35 <sub>0.11</sub>	17.65 <sub>7.90</sub>	14.21 <sub>4.88</sub>	17.31 <sub>0.55</sub>	22.03 <sub>1.67</sub>	17.86 <sub>2.16</sub>	17.87 <sub>2.35</sub>	5.43 <sub>1.65</sub>	14.39 <sub>0.67</sub>	6.90 <sub>0.26</sub>	28.80 <sub>1.95</sub>
	Test	w/o	52.61 <sub>0.00</sub>	51.74 <sub>0.00</sub>	32.92 <sub>0.00</sub>	14.46 <sub>0.00</sub>	19.68 <sub>0.00</sub>	27.22 <sub>0.00</sub>	39.90 <sub>0.00</sub>	10.44 <sub>0.00</sub>	8.01 <sub>0.00</sub>	22.48 <sub>0.00</sub>	12.02 <sub>0.00</sub>	49.70 <sub>0.00</sub>
GPT-2	MF1 (%)	Val.	w/o	57.49 <sub>2.32</sub>	54.24 <sub>0.30</sub>	17.41 <sub>10.77</sub>	20.08 <sub>5.50</sub>	20.17 <sub>0.46</sub>	25.34 <sub>0.77</sub>	21.52 <sub>3.83</sub>	8.81 <sub>0.48</sub>	3.94 <sub>2.29</sub>	11.62 <sub>0.67</sub>	12.48 <sub>0.25</sub>
		w/	54.24 <sub>0.30</sub>	17.41 <sub>10.77</sub>	20.08 <sub>5.50</sub>	20.17 <sub>0.46</sub>	25.34 <sub>0.77</sub>	21.52 <sub>3.83</sub>	8.81 <sub>0.48</sub>	3.94 <sub>2.29</sub>	11.62 <sub>0.67</sub>	12.48 <sub>0.25</sub>	14.84 <sub>2.55</sub>	22.33

In terms of implementation, we use the `MulticlassCalibrationError` module given by the TorchMetrics <sup>3</sup> with `n_bins=10, norm='l1'`.

## G Complete Results: Reliability Diagrams and Confident Distribution

**Reliability diagrams:** as shown in Fig. 9 - 14 for GPT-J, and Fig. 21 - 24 for GPT-2.

**Confidence distributions:** as shown in Fig. 15 - 20 for GPT-J, and Fig. 25 - 28 for GPT-2.

## H Case Analysis: NOISYICL Furtherance Correct ICL

Moreover, we find that in some cases, unperturbed ICL can't benefit correctly from scaling the number of demos, while, NOISYICL can help the model correct this issue, as shown in Fig. 5. These unperturbed models exhibit an overfitting-like phenomenon and also low accuracies, while NOISYICL can relieve it.

We speculate the reason is the mismatch between the pre-trained knowledge and ICL inputs. This leads to a decrease in the model's in-context task learning (Pan et al., 2023) ability, while NOISYICL reduces such a gap between pre-trained data and ICL style data, which makes models extract information from ICL inputs better.

<sup>3</sup>[lightning.ai/docs/torchmetrics/stable/classification/calibration\\_error.html#multiclasscalibrationerror](https://lightning.ai/docs/torchmetrics/stable/classification/calibration_error.html#multiclasscalibrationerror)

## I License for Artifacts

Here we discuss the license of the artifacts used in this paper.

**Models.** GPT-2 is under the MIT license, and GPT-J is under the apache-2.0 license.

**Datasets.** We list the open-source license for the datasets used in this paper as follows:

- cc-by-4.0: PS, SE'14R, SE'14L, TEE, TES, TEH
- cc-by-sa-3.0: HS, SST2, RTE, MRPC, FP
- agpl-3.0: Ethos

**Consistency of Usage.** Models are used with their original usage. Due to the different data splitting of these datasets, to ensure the consistency of the experiment methods, we use a re-splitting method as described in Appendix A. However, the overall usage is consistent with their intended use.

Table 4: Accuracy and Macro-F1 results (% ,  $mean_{std}$ ,  $k = 1$ ). A better result is in green. Notation is the same with the Tabel 1.

Dataset		PS	HS	SE'14R	SE'14L	RTE	MRPC	Ethos	FP	SST2	TEE	TES	TEH	Mean
	$\lambda$	0.2	0.012	0.1	0.02	0.2	0.016	0.016	0.014	0.018	0.014	0.002	0.08	—
Acc. (%)	Val. w/o	39.43 <sub>0.99</sub>	76.02 <sub>0.86</sub>	57.85 <sub>1.45</sub>	53.77 <sub>1.82</sub>	50.41 <sub>1.15</sub>	58.48 <sub>1.48</sub>	50.32 <sub>1.81</sub>	53.93 <sub>0.82</sub>	68.20 <sub>0.52</sub>	43.48 <sub>0.80</sub>	43.83 <sub>2.45</sub>	50.90 <sub>1.02</sub>	52.22
	w/	53.24 <sub>4.76</sub>	76.86 <sub>0.54</sub>	64.28 <sub>2.15</sub>	55.57 <sub>1.05</sub>	51.76 <sub>1.55</sub>	58.28 <sub>1.35</sub>	51.00 <sub>1.77</sub>	53.96 <sub>1.27</sub>	71.56 <sub>1.07</sub>	44.32 <sub>0.50</sub>	46.37 <sub>1.09</sub>	52.13 <sub>0.96</sub>	56.61
	Test w/o	35.08 <sub>1.74</sub>	84.92 <sub>0.94</sub>	53.66 <sub>1.12</sub>	51.60 <sub>1.29</sub>	50.52 <sub>1.31</sub>	53.50 <sub>2.26</sub>	50.40 <sub>1.37</sub>	52.29 <sub>1.48</sub>	73.89 <sub>0.84</sub>	41.94 <sub>0.92</sub>	44.80 <sub>1.26</sub>	50.71 <sub>1.46</sub>	53.61
MF1 (%)	Val. w/o	26.01 <sub>0.70</sub>	24.73 <sub>0.66</sub>	50.78 <sub>1.61</sub>	52.26 <sub>1.85</sub>	50.36 <sub>1.18</sub>	51.11 <sub>1.43</sub>	49.04 <sub>1.83</sub>	38.25 <sub>1.01</sub>	67.39 <sub>0.58</sub>	26.71 <sub>1.08</sub>	38.75 <sub>2.41</sub>	49.42 <sub>0.85</sub>	43.73
	w/	23.18 <sub>1.34</sub>	24.49 <sub>0.92</sub>	52.49 <sub>2.80</sub>	53.61 <sub>1.35</sub>	48.28 <sub>1.74</sub>	50.61 <sub>1.72</sub>	49.88 <sub>1.81</sub>	38.94 <sub>1.99</sub>	71.02 <sub>1.16</sub>	27.65 <sub>1.31</sub>	41.43 <sub>1.61</sub>	50.75 <sub>0.97</sub>	44.36
	Test w/o	25.30 <sub>1.20</sub>	24.73 <sub>0.71</sub>	47.91 <sub>1.24</sub>	50.62 <sub>1.25</sub>	50.17 <sub>1.19</sub>	49.97 <sub>2.33</sub>	49.81 <sub>1.33</sub>	37.90 <sub>1.98</sub>	72.25 <sub>1.04</sub>	24.67 <sub>1.43</sub>	40.10 <sub>1.28</sub>	48.45 <sub>1.61</sub>	43.49
ECE <sub>1</sub> (%)	Val. w/o	23.09 <sub>1.25</sub>	22.67 <sub>0.82</sub>	11.25 <sub>1.73</sub>	17.01 <sub>1.66</sub>	49.33 <sub>1.14</sub>	35.31 <sub>1.39</sub>	47.42 <sub>1.82</sub>	27.54 <sub>0.97</sub>	5.33 <sub>1.23</sub>	44.33 <sub>1.12</sub>	28.71 <sub>2.16</sub>	45.64 <sub>1.19</sub>	29.80
	w/	16.64 <sub>1.70</sub>	21.81 <sub>0.59</sub>	13.33 <sub>2.62</sub>	15.57 <sub>0.64</sub>	26.82 <sub>5.71</sub>	35.56 <sub>1.89</sub>	46.63 <sub>1.82</sub>	27.39 <sub>1.11</sub>	3.72 <sub>1.19</sub>	42.88 <sub>0.69</sub>	26.36 <sub>1.40</sub>	43.95 <sub>1.10</sub>	26.72
	Test w/o	28.18 <sub>1.95</sub>	13.97 <sub>0.95</sub>	12.05 <sub>1.52</sub>	15.95 <sub>1.05</sub>	49.20 <sub>1.34</sub>	39.94 <sub>2.33</sub>	47.35 <sub>1.40</sub>	28.24 <sub>1.42</sub>	4.60 <sub>0.95</sub>	45.34 <sub>0.90</sub>	27.63 <sub>0.85</sub>	46.98 <sub>1.37</sub>	29.95
GPT-J	w/	20.85 <sub>2.58</sub>	13.21 <sub>1.07</sub>	15.45 <sub>1.84</sub>	14.28 <sub>0.93</sub>	35.64 <sub>3.69</sub>	35.56 <sub>1.89</sub>	46.76 <sub>2.10</sub>	28.13 <sub>0.97</sub>	5.53 <sub>1.44</sub>	44.71 <sub>0.93</sub>	27.08 <sub>1.93</sub>	45.28 <sub>1.33</sub>	26.87
	$\lambda$	0.02	0.002	0.016	0.004	0.002	0.002	0.2	0.06	0.02	0.02	0.06	0.1	—
	Val. w/o	42.42 <sub>0.97</sub>	41.37 <sub>1.11</sub>	45.68 <sub>1.09</sub>	39.69 <sub>1.69</sub>	50.62 <sub>1.35</sub>	68.71 <sub>0.20</sub>	47.26 <sub>0.72</sub>	45.45 <sub>1.81</sub>	57.79 <sub>1.23</sub>	28.11 <sub>0.82</sub>	41.45 <sub>1.36</sub>	42.70 <sub>0.32</sub>	45.94
Acc. (%)	w/	45.55 <sub>3.18</sub>	37.34 <sub>1.33</sub>	48.75 <sub>1.80</sub>	41.66 <sub>1.04</sub>	50.18 <sub>2.11</sub>	68.81 <sub>0.18</sub>	53.35 <sub>2.99</sub>	55.33 <sub>2.72</sub>	59.26 <sub>1.38</sub>	32.73 <sub>0.53</sub>	45.25 <sub>1.34</sub>	48.41 <sub>2.15</sub>	48.88
	Test w/o	35.07 <sub>1.44</sub>	39.64 <sub>0.61</sub>	43.24 <sub>1.43</sub>	38.47 <sub>0.97</sub>	50.22 <sub>1.27</sub>	62.94 <sub>0.42</sub>	48.65 <sub>0.59</sub>	45.09 <sub>1.51</sub>	59.76 <sub>1.67</sub>	27.61 <sub>1.20</sub>	37.98 <sub>1.24</sub>	44.37 <sub>0.81</sub>	44.42
	w/	40.08 <sub>2.90</sub>	34.92 <sub>1.12</sub>	44.36 <sub>0.95</sub>	39.41 <sub>1.30</sub>	50.58 <sub>2.31</sub>	62.47 <sub>0.29</sub>	51.69 <sub>2.95</sub>	52.12 <sub>1.60</sub>	60.64 <sub>1.78</sub>	30.55 <sub>0.79</sub>	43.88 <sub>1.53</sub>	55.35 <sub>0.79</sub>	47.17
MF1 (%)	Val. w/o	26.44 <sub>0.73</sub>	20.32 <sub>1.44</sub>	32.96 <sub>1.92</sub>	33.69 <sub>1.92</sub>	50.59 <sub>1.35</sub>	41.67 <sub>0.39</sub>	43.18 <sub>0.82</sub>	30.19 <sub>1.51</sub>	56.71 <sub>1.25</sub>	22.59 <sub>0.81</sub>	36.97 <sub>1.57</sub>	32.23 <sub>0.53</sub>	35.63
	w/	26.37 <sub>1.41</sub>	18.95 <sub>1.29</sub>	35.97 <sub>1.24</sub>	35.61 <sub>1.25</sub>	50.13 <sub>2.10</sub>	42.28 <sub>0.48</sub>	47.26 <sub>2.81</sub>	34.47 <sub>0.77</sub>	59.22 <sub>1.39</sub>	22.87 <sub>1.07</sub>	30.49 <sub>1.97</sub>	51.68 <sub>2.03</sub>	37.94
	Test w/o	23.97 <sub>1.23</sub>	17.45 <sub>0.85</sub>	33.27 <sub>1.93</sub>	36.32 <sub>1.03</sub>	49.85 <sub>1.26</sub>	39.41 <sub>0.66</sub>	43.80 <sub>0.77</sub>	30.92 <sub>1.38</sub>	59.63 <sub>1.61</sub>	22.32 <sub>1.38</sub>	36.21 <sub>1.30</sub>	38.79 <sub>1.06</sub>	36.00
GPT-2	w/	26.31 <sub>1.55</sub>	15.87 <sub>1.14</sub>	33.78 <sub>1.35</sub>	37.11 <sub>1.53</sub>	50.18 <sub>2.40</sub>	39.10 <sub>0.41</sub>	44.93 <sub>3.66</sub>	34.58 <sub>1.06</sub>	59.52 <sub>1.82</sub>	21.69 <sub>1.08</sub>	33.40 <sub>1.36</sub>	50.57 <sub>1.50</sub>	37.25
	Val. w/o	14.45 <sub>1.49</sub>	25.96 <sub>1.14</sub>	16.68 <sub>1.33</sub>	25.55 <sub>1.99</sub>	42.92 <sub>1.41</sub>	19.22 <sub>0.49</sub>	35.57 <sub>1.01</sub>	21.11 <sub>1.76</sub>	5.40 <sub>1.16</sub>	31.62 <sub>1.17</sub>	20.72 <sub>1.70</sub>	46.15 <sub>0.38</sub>	25.44
	w/	9.64 <sub>1.86</sub>	30.37 <sub>1.36</sub>	17.27 <sub>2.01</sub>	23.65 <sub>0.78</sub>	43.31 <sub>2.03</sub>	17.87 <sub>0.49</sub>	19.53 <sub>3.80</sub>	14.80 <sub>0.62</sub>	3.94 <sub>0.41</sub>	30.66 <sub>0.71</sub>	24.11 <sub>2.54</sub>	26.44 <sub>2.47</sub>	21.80
ECE <sub>1</sub> (%)	Test w/o	18.82 <sub>1.59</sub>	26.53 <sub>0.96</sub>	18.64 <sub>1.41</sub>	24.56 <sub>0.96</sub>	43.56 <sub>1.28</sub>	24.09 <sub>0.58</sub>	35.15 <sub>0.76</sub>	22.13 <sub>1.42</sub>	2.83 <sub>1.04</sub>	32.31 <sub>1.34</sub>	22.45 <sub>1.38</sub>	39.96 <sub>0.77</sub>	25.92
	w/	13.20 <sub>2.03</sub>	32.02 <sub>1.47</sub>	20.45 <sub>1.07</sub>	23.83 <sub>1.22</sub>	43.07 <sub>2.30</sub>	23.01 <sub>0.87</sub>	23.65 <sub>6.52</sub>	14.48 <sub>1.15</sub>	3.22 <sub>1.50</sub>	33.43 <sub>1.60</sub>	19.37 <sub>1.64</sub>	23.89 <sub>1.50</sub>	22.80

Table 5: Accuracy and Macro-F1 results (% ,  $mean_{std}$ ,  $k = 2$ ). A better result is in green. Notation is the same with the Tabel 1.

Dataset		PS	HS	SE'14R	SE'14L	RTE	MRPC	Ethos	FP	SST2	TEE	TES	TEH	Mean
	$\lambda$	0.2	0.2	0.1	0.1	0.08	0.2	0.014	0.08	0.016	0.018	0.02	0.008	—
Acc. (%)	Val. w/o	42.54 <sub>1.49</sub>	72.03 <sub>0.66</sub>	41.50 <sub>0.97</sub>	40.88 <sub>0.49</sub>	51.31 <sub>1.19</sub>	55.74 <sub>1.24</sub>	49.81 <sub>2.28</sub>	61.84 <sub>0.60</sub>	72.05 <sub>0.96</sub>	44.47 <sub>0.62</sub>	48.85 <sub>0.75</sub>	49.08 <sub>1.16</sub>	52.51
	w/	57.11 <sub>4.74</sub>	77.34 <sub>3.29</sub>	58.57 <sub>3.39</sub>	51.25 <sub>2.71</sub>	50.76 <sub>0.99</sub>	58.98 <sub>6.13</sub>	51.37 <sub>1.15</sub>	60.47 <sub>1.77</sub>	73.38 <sub>0.70</sub>	44.67 <sub>0.63</sub>	49.14 <sub>0.78</sub>	50.62 <sub>2.11</sub>	56.97
	Test w/o	39.23 <sub>0.98</sub>	82.42 <sub>1.21</sub>	36.69 <sub>1.30</sub>	39.26 <sub>1.67</sub>	48.13 <sub>2.60</sub>	50.97 <sub>1.43</sub>	51.41 <sub>1.93</sub>	60.77 <sub>0.64</sub>	72.60 <sub>1.66</sub>	44.28 <sub>0.73</sub>	48.79 <sub>1.00</sub>	51.68 <sub>1.95</sub>	52.19
MF1 (%)	Val. w/o	25.17 <sub>1.10</sub>	25.03 <sub>0.33</sub>	37.42 <sub>1.13</sub>	39.11 <sub>0.83</sub>	50.44 <sub>1.23</sub>	51.35 <sub>1.45</sub>	49.78 <sub>2.28</sub>	31.83 <sub>1.31</sub>	71.32 <sub>1.07</sub>	25.43 <sub>1.03</sub>	36.69 <sub>1.16</sub>	49.05 <sub>1.17</sub>	41.05
	w/	22.75 <sub>1.84</sub>	24.15 <sub>0.56</sub>	50.06 <sub>1.11</sub>	49.26 <sub>2.50</sub>	50.47 <sub>2.22</sub>	48.95 <sub>0.99</sub>	51.36 <sub>1.15</sub>	35.39 <sub>3.22</sub>	72.77 <sub>0.78</sub>	25.16 <sub>0.84</sub>	37.72 <sub>1.56</sub>	50.56 <sub>2.12</sub>	43.22
	Test w/o	25.51 <sub>1.25</sub>	24.95 <sub>0.85</sub>	34.99 <sub>1.38</sub>	39.75 <sub>1.63</sub>	47.82 <sub>2.64</sub>	49.56 <sub>1.48</sub>	51.18 <sub>1.96</sub>	30.66 <sub>0.98</sub>	70.34 <sub>1.87</sub>	24.05 <sub>0.97</sub>	36.78 <sub>1.28</sub>	51.49 <sub>1.96</sub>	40.59
ECE <sub>1</sub> (%)	Val. w/o	15.11 <sub>29</sub>	21.52 <sub>0.58</sub>	23.45 <sub>1.46</sub>	27.07 <sub>0.89</sub>	39.26 <sub>0.83</sub>	23.82 <sub>1.59</sub>	38.51 <sub>2.24</sub>	22.54 <sub>0.77</sub>	3.56 <sub>0.89</sub>	44.33 <sub>0.66</sub>	19.22 <sub>0.62</sub>	40.33 <sub>1.30</sub>	26.56
	w/	13.69 <sub>3.31</sub>	11.25 <sub>5.10</sub>	8.03 <sub>4.49</sub>	15.54 <sub>3.79</sub>	39.51 <sub>3.65</sub>	12.12 <sub>2.82</sub>	37.12 <sub>0.98</sub>	21.10 <sub>1.70</sub>	4.80 <sub>0.74</sub>	43.87 <sub>0.96</sub>	18.64 <sub>1.08</sub>	38.66 <sub>2.22</sub>	22.03
	Test w/o	18.59 <sub>1.01</sub>	13.15 <sub>1.15</sub>	27.87 <sub>1.35</sub>	26.71 <sub>1.83</sub>	41.77 <sub>2.91</sub>	28.29 <sub>1.47</sub>	36.62 <sub>1.76</sub>	22.97 <sub>0.87</sub>	4.77 <sub>0.84</sub>	44.53 <sub>0.98</sub>	19.92 <sub>1.25</sub>	36.63 <sub>1.86</sub>	26.82
GPT-J	w/	13.49 <sub>2.88</sub>	6.52 <sub>2.25</sub>	11.44 <sub>3.91</sub>	15.80 <sub>2.76</sub>	41.33 <sub>2.22</sub>	20.24 <sub>6.56</sub>	36.61 <sub>1.67</sub>	21.88 <sub>1.75</sub>	5.32 <sub>1.16</sub>	44.40 <sub>0.77</sub>	19.98 <sub>0.70</sub>	37.48 <sub>1.92</sub>	22.87
	$\lambda$	0.02	0.1	0.08	0.08	0.006	0.002	0.2	0.04	0.02	0.02	0.04	0.2	—
	Val. w/o	41.86 <sub>1.75</sub>	48.24 <sub>1.37</sub>	46.99 <sub>0.84</sub>	42.42 <sub>1.10</sub>	50.88 <sub>1.05</sub>	68.55 <sub>0.16</sub>	45.85 <sub>1.27</sub>	49.10 <sub>1.15</sub>	56.56 <sub>1.31</sub>	29.06 <sub>0.71</sub>	39.39 <sub>1.51</sub>	42.81 <sub>0.63</sub>	46.81
Acc. (%)	w/	47.23 <sub>3.35</sub>	48.75 <sub>6.40</sub>	52.85 <sub>8.89</sub>	46.62 <sub>1.59</sub>	51.80 <sub>1.48</sub>	68.89 <sub>0.43</sub>	52.84 <sub>2.43</sub>	61.31 <sub>0.81</sub>	59.45 <sub>1.20</sub>	30.29 <sub>0.82</sub>	45.70 <sub>1.46</sub>	51.09 <sub>4.07</sub>	51.40
	Test w/o	38.26 <sub>1.37</sub>	52.04 <sub>0.91</sub>	42.52 <sub>1.18</sub>	39.71 <sub>1.47</sub>	49.67 <sub>0.80</sub>	62.77 <sub>0.39</sub>	48.26 <sub>1.00</sub>	47.42 <sub>1.45</sub>	58.96 <sub>1.07</sub>	28.47 <sub>0.97</sub>	32.11 <sub>1.37</sub>	44.74 <sub>0.82</sub>	45.41
	w/	40.95 <sub>1.84</sub>	45.52 <sub>3.99</sub>	50.12 <sub>3.58</sub>	40.31 <sub>0.90</sub>	50.64 <sub>1.87</sub>	57.34 <sub>1.01</sub>	51.76 <sub>2.27</sub>	59.41 <sub>1.11</sub>	61.33 <sub>1.51</sub>	28.49 <sub>1.04</sub>	40.29 <sub>3.18</sub>	53.01 <sub>3.31</sub>	48.26
MF1 (%)	Val. w/o	25.20 <sub>1.70</sub>	21.61 <sub>1.20</sub>	36.33 <sub>1.10</sub>	36.50 <sub>0.51</sub>	50.38 <sub>1.03</sub>	41.38 <sub>0.39</sub>	41.44 <sub>1.75</sub>						

Table 6: Accuracy and Macro-F1 results (%),  $mean_{std}$ ,  $k = 8$ . A better result is in green. Notation is the same with the Tabel 1. Some experiments can't be conducted due to the length of the input sequence.

Dataset		PS	HS	SE'14R	SE'14L	RTE	MRPC	Ethos	FP	SST2	TEE	TES	TEH	Mean
$\lambda$		0.2	0.2	0.1	0.1	0.1	0.2	0.004	0.012	0.002	0.006	0.08	0.002	—
Acc. (%)	Val. w/o	27.32 <sub>0.84</sub>	68.87 <sub>1.30</sub>	35.61 <sub>0.57</sub>	33.57 <sub>0.74</sub>	50.12 <sub>0.29</sub>	35.86 <sub>0.55</sub>	57.29 <sub>1.69</sub>	62.34 <sub>0.05</sub>	83.96 <sub>0.86</sub>	47.64 <sub>0.52</sub>	51.13 <sub>0.66</sub>	50.62 <sub>1.62</sub>	50.37
	Val. w/	62.01 <sub>4.95</sub>	85.29 <sub>0.95</sub>	45.84 <sub>1.13</sub>	40.84 <sub>6.41</sub>	50.49 <sub>0.89</sub>	58.87 <sub>7.93</sub>	57.95 <sub>1.35</sub>	62.44 <sub>0.08</sub>	83.89 <sub>0.79</sub>	47.81 <sub>0.90</sub>	51.35 <sub>1.67</sub>	52.05 <sub>0.75</sub>	58.24
	Test w/o	25.70 <sub>0.84</sub>	77.10 <sub>1.75</sub>	33.30 <sub>0.63</sub>	30.80 <sub>0.71</sub>	46.35 <sub>1.39</sub>	39.35 <sub>0.84</sub>	61.67 <sub>1.38</sub>	61.76 <sub>0.06</sub>	80.59 <sub>1.02</sub>	46.37 <sub>0.50</sub>	50.68 <sub>0.35</sub>	54.57 <sub>1.41</sub>	50.69
GPT-J MF1 (%)	Val. w/o	16.07 <sub>0.47</sub>	23.66 <sub>0.32</sub>	30.15 <sub>1.14</sub>	24.48 <sub>1.40</sub>	45.66 <sub>0.49</sub>	33.37 <sub>0.56</sub>	56.68 <sub>1.81</sub>	25.91 <sub>0.26</sub>	83.80 <sub>0.91</sub>	27.20 <sub>0.59</sub>	33.72 <sub>1.47</sub>	48.52 <sub>1.98</sub>	37.44
	Val. w/	20.91 <sub>1.79</sub>	27.76 <sub>1.35</sub>	41.70 <sub>10.30</sub>	35.28 <sub>8.07</sub>	45.89 <sub>3.44</sub>	46.69 <sub>2.90</sub>	57.45 <sub>1.37</sub>	26.24 <sub>0.36</sub>	83.72 <sub>0.82</sub>	27.39 <sub>0.84</sub>	36.84 <sub>5.55</sub>	49.94 <sub>0.99</sub>	41.65
	Test w/o	16.01 <sub>0.94</sub>	26.26 <sub>1.29</sub>	28.72 <sub>0.67</sub>	27.21 <sub>0.78</sub>	42.58 <sub>1.37</sub>	33.58 <sub>0.98</sub>	60.88 <sub>1.44</sub>	25.60 <sub>0.25</sub>	78.98 <sub>1.18</sub>	24.50 <sub>0.70</sub>	33.19 <sub>0.40</sub>	54.36 <sub>1.42</sub>	37.66
ECE <sub>1</sub> (%)	Val. w/o	42.98 <sub>0.85</sub>	10.43 <sub>1.20</sub>	32.79 <sub>0.53</sub>	42.43 <sub>0.99</sub>	21.59 <sub>0.34</sub>	33.19 <sub>0.72</sub>	13.63 <sub>1.44</sub>	24.54 <sub>0.13</sub>	12.61 <sub>0.98</sub>	40.10 <sub>0.65</sub>	13.99 <sub>0.17</sub>	21.11 <sub>1.51</sub>	25.78
	Val. w/	11.81 <sub>4.56</sub>	9.51 <sub>2.14</sub>	19.27 <sub>11.45</sub>	25.96 <sub>10.94</sub>	23.13 <sub>2.74</sub>	8.10 <sub>2.84</sub>	12.47 <sub>1.21</sub>	24.33 <sub>0.69</sub>	12.80 <sub>0.75</sub>	39.95 <sub>1.00</sub>	7.25 <sub>4.79</sub>	20.01 <sub>0.66</sub>	17.88
	Test w/o	44.92 <sub>0.80</sub>	7.93 <sub>1.41</sub>	36.15 <sub>0.98</sub>	43.80 <sub>1.01</sub>	25.13 <sub>1.62</sub>	29.83 <sub>0.74</sub>	9.35 <sub>1.45</sub>	25.00 <sub>0.23</sub>	10.52 <sub>0.94</sub>	41.83 <sub>0.54</sub>	14.13 <sub>0.43</sub>	16.78 <sub>1.45</sub>	25.45
GPT-J ECE <sub>1</sub> (%)	Test w/	8.86 <sub>3.57</sub>	6.04 <sub>1.94</sub>	22.26 <sub>9.50</sub>	29.44 <sub>8.54</sub>	28.55 <sub>6.11</sub>	15.28 <sub>7.58</sub>	9.77 <sub>0.86</sub>	24.92 <sub>0.52</sub>	9.88 <sub>0.66</sub>	41.83 <sub>0.80</sub>	9.24 <sub>2.47</sub>	16.77 <sub>0.95</sub>	18.57

Table 7: Accuracy and Macro-F1 results (%),  $mean_{std}$ ,  $k = 16$ . A better result is in green. Notation is the same with the Tabel 1. Some experiments can't be conducted due to the length of the input sequence.

Dataset		PS	HS	SE'14R	SE'14L	RTE	FP	SST2	TEE	TES	TEH	Mean
$\lambda$		0.2	0.2	0.1	0.1	0.2	0.1	0.004	0.008	0.08	0.002	—
Acc. (%)	Val. w/o	18.52 <sub>0.78</sub>	70.02 <sub>0.52</sub>	40.86 <sub>0.90</sub>	34.59 <sub>0.41</sub>	48.91 <sub>0.62</sub>	62.42 <sub>0.10</sub>	88.24 <sub>0.70</sub>	49.22 <sub>0.47</sub>	51.84 <sub>0.97</sub>	54.59 <sub>0.87</sub>	51.92
	Val. w/	61.78 <sub>5.94</sub>	85.98 <sub>0.61</sub>	46.84 <sub>12.92</sub>	43.01 <sub>8.67</sub>	50.76 <sub>1.15</sub>	62.66 <sub>0.35</sub>	88.83 <sub>0.71</sub>	49.30 <sub>0.53</sub>	53.11 <sub>1.99</sub>	55.20 <sub>0.91</sub>	59.75
	Test w/o	18.75 <sub>0.38</sub>	74.97 <sub>1.43</sub>	34.98 <sub>1.10</sub>	30.00 <sub>0.58</sub>	45.95 <sub>0.45</sub>	61.82 <sub>0.09</sub>	81.88 <sub>1.00</sub>	47.11 <sub>0.35</sub>	51.70 <sub>0.42</sub>	55.90 <sub>0.81</sub>	50.31
GPT-J MF1 (%)	Val. w/o	10.33 <sub>0.46</sub>	30.44 <sub>0.84</sub>	32.79 <sub>0.78</sub>	24.15 <sub>0.66</sub>	38.73 <sub>0.71</sub>	26.08 <sub>0.46</sub>	88.22 <sub>0.70</sub>	28.75 <sub>0.40</sub>	36.01 <sub>1.76</sub>	53.61 <sub>0.88</sub>	36.91
	Val. w/	20.00 <sub>0.71</sub>	23.54 <sub>0.57</sub>	41.11 <sub>11.38</sub>	35.07 <sub>10.38</sub>	42.85 <sub>7.50</sub>	27.18 <sub>1.53</sub>	88.81 <sub>0.72</sub>	29.04 <sub>0.55</sub>	40.22 <sub>6.60</sub>	54.02 <sub>1.02</sub>	40.18
	Test w/o	9.61 <sub>0.26</sub>	27.64 <sub>1.69</sub>	28.80 <sub>0.97</sub>	25.05 <sub>0.94</sub>	35.57 <sub>0.69</sub>	25.80 <sub>0.31</sub>	80.33 <sub>1.12</sub>	25.50 <sub>0.54</sub>	35.44 <sub>0.90</sub>	55.76 <sub>0.81</sub>	34.95
GPT-J ECE <sub>1</sub> (%)	Val. w/o	54.29 <sub>0.95</sub>	4.01 <sub>0.77</sub>	27.53 <sub>1.27</sub>	40.90 <sub>0.52</sub>	26.13 <sub>0.54</sub>	26.34 <sub>0.14</sub>	14.43 <sub>0.75</sub>	38.97 <sub>0.58</sub>	11.46 <sub>0.99</sub>	13.66 <sub>0.72</sub>	25.77
	Val. w/	11.75 <sub>5.60</sub>	9.70 <sub>2.72</sub>	18.56 <sub>10.81</sub>	23.33 <sub>13.20</sub>	21.01 <sub>13.24</sub>	18.07 <sub>5.22</sub>	15.49 <sub>0.58</sub>	38.61 <sub>0.69</sub>	6.12 <sub>2.85</sub>	13.25 <sub>0.68</sub>	17.59
	Test w/o	53.72 <sub>0.37</sub>	4.51 <sub>0.87</sub>	34.09 <sub>1.27</sub>	43.06 <sub>0.62</sub>	29.93 <sub>0.68</sub>	26.90 <sub>0.19</sub>	9.96 <sub>1.00</sub>	40.84 <sub>0.45</sub>	11.24 <sub>0.41</sub>	13.45 <sub>0.86</sub>	26.77
	Test w/	8.33 <sub>4.47</sub>	5.63 <sub>2.27</sub>	20.29 <sub>12.35</sub>	24.94 <sub>10.05</sub>	17.53 <sub>7.30</sub>	20.44 <sub>3.90</sub>	9.99 <sub>0.89</sub>	41.00 <sub>0.83</sub>	5.19 <sub>2.52</sub>	13.03 <sub>1.33</sub>	16.64

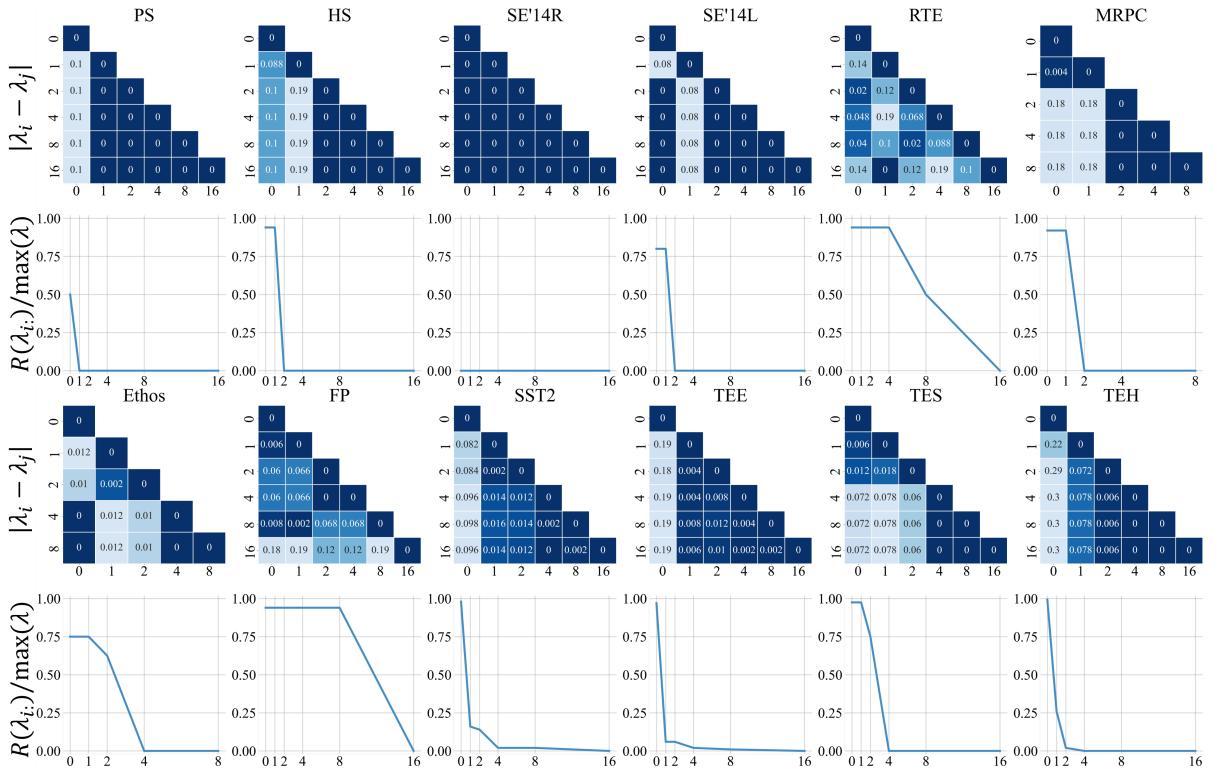


Figure 7: The stability of the optimal  $\lambda$  on various  $k$  for GPT-J. Upper figure: the distance between the optimal  $\lambda$  of different demo numbers; Lower figure: the normalized remaining range.

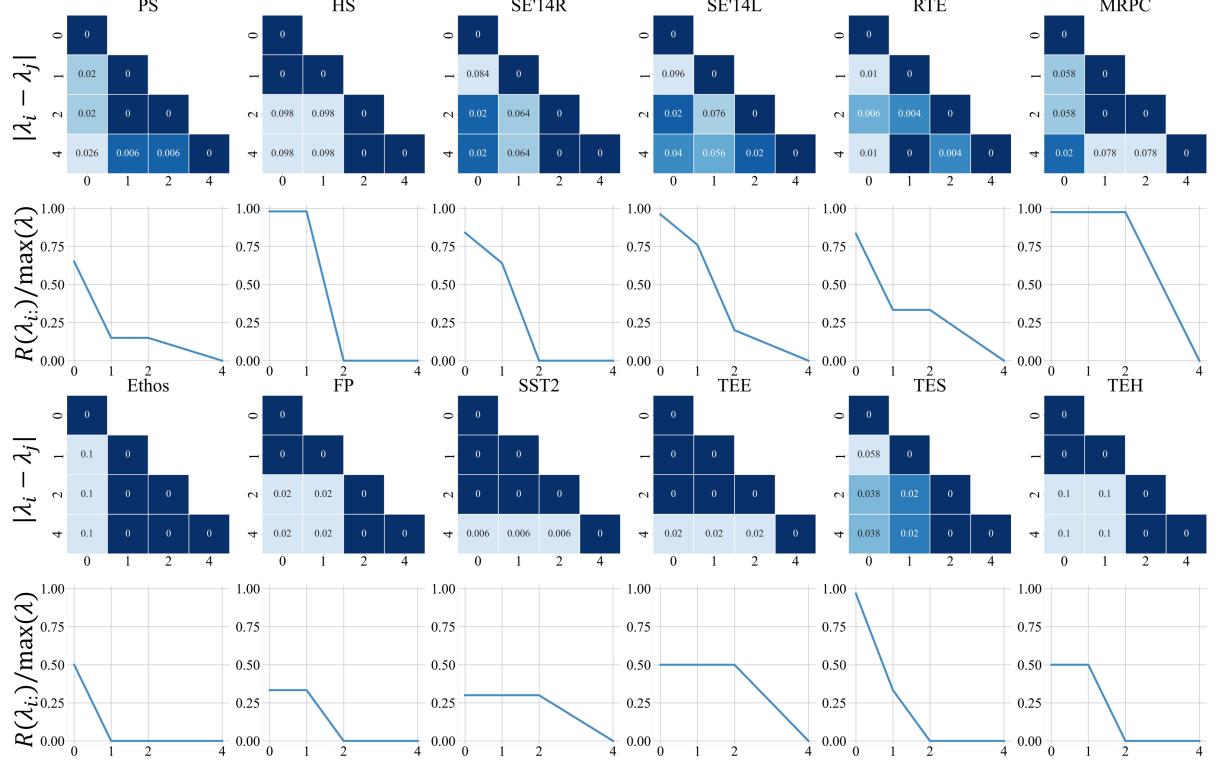


Figure 8: The stability of the optimal  $\lambda$  on various  $k$  for GPT-2. Upper figure: the distance between the optimal  $\lambda$  of different demo numbers; Lower figure: the normalized remaining range.

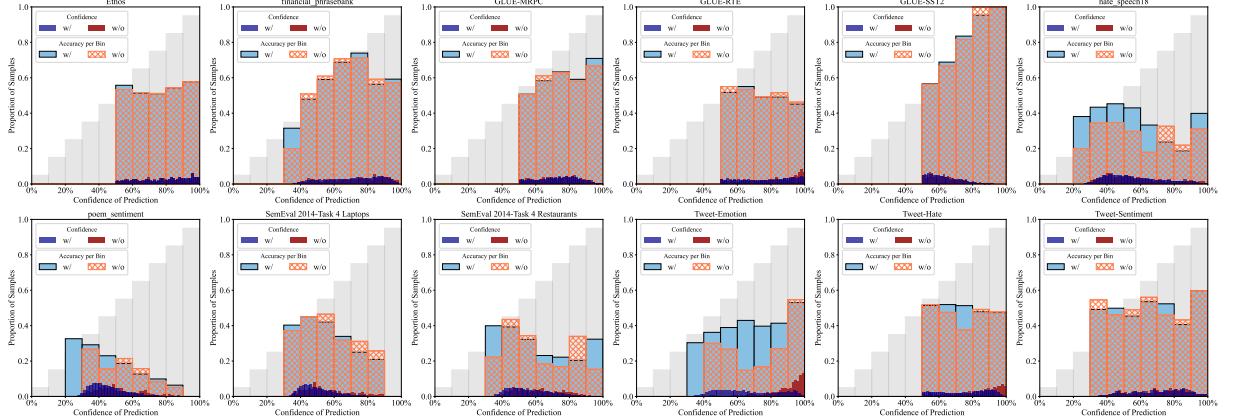


Figure 9: The reliability diagrams of GPT-J ( $k = 0$ ).

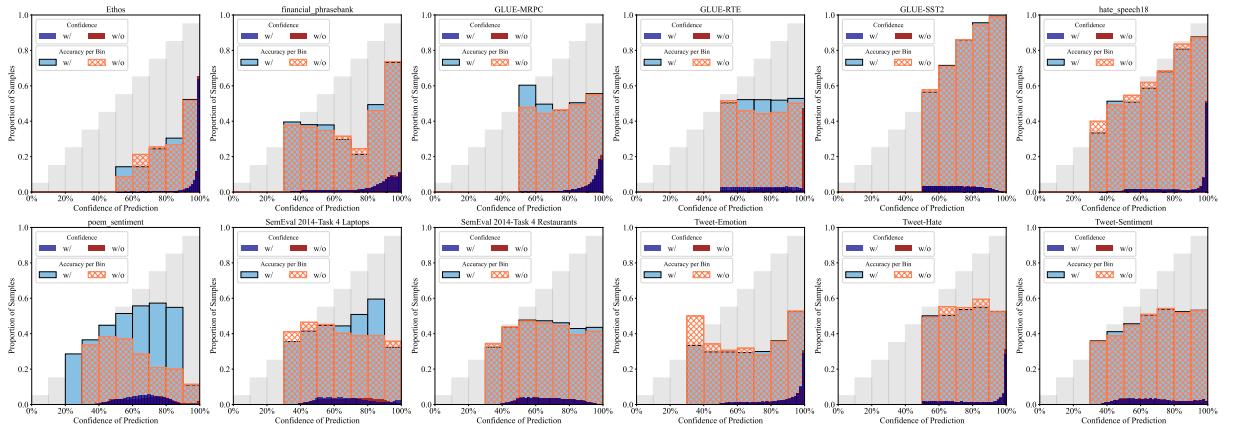


Figure 10: The reliability diagrams of GPT-J ( $k = 1$ ).

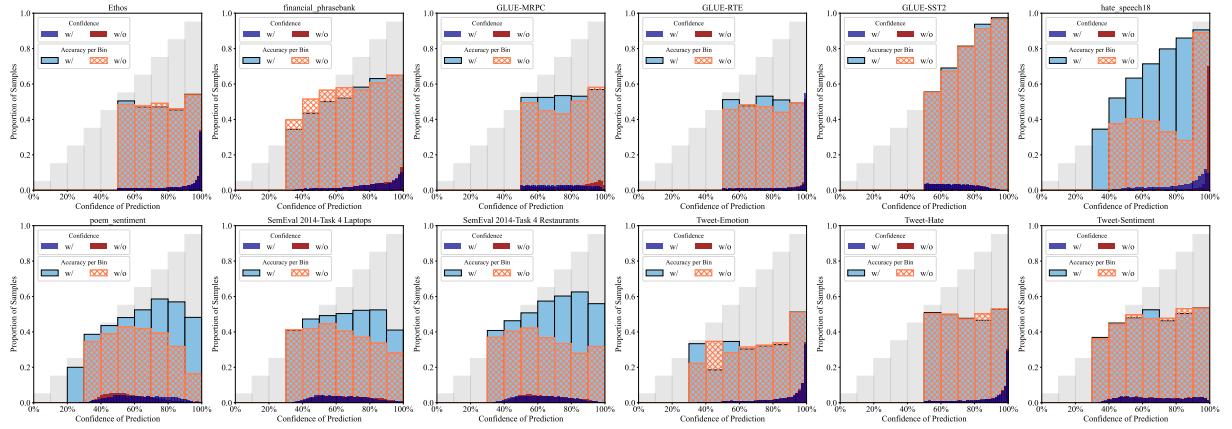


Figure 11: The reliability diagrams of GPT-J ( $k = 2$ ).

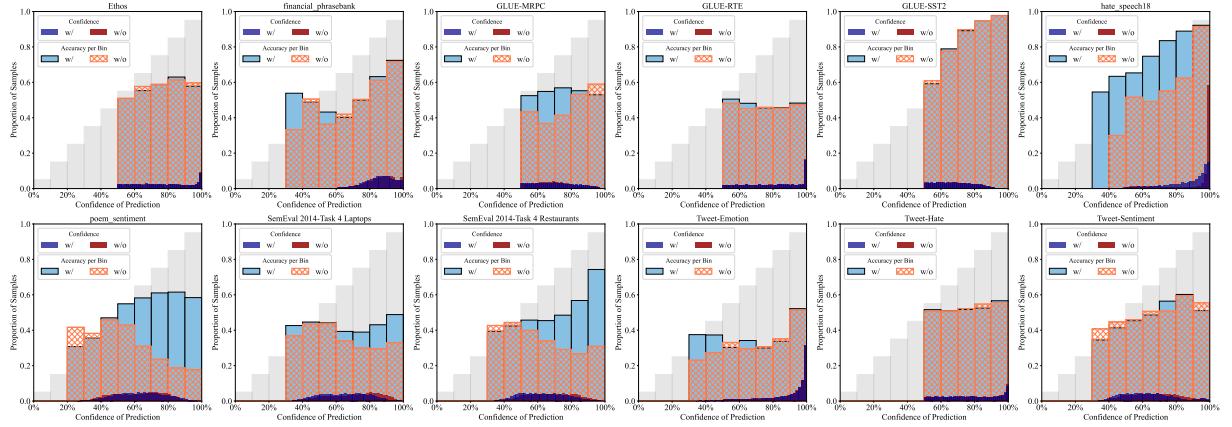


Figure 12: The reliability diagrams of GPT-J ( $k = 4$ ).

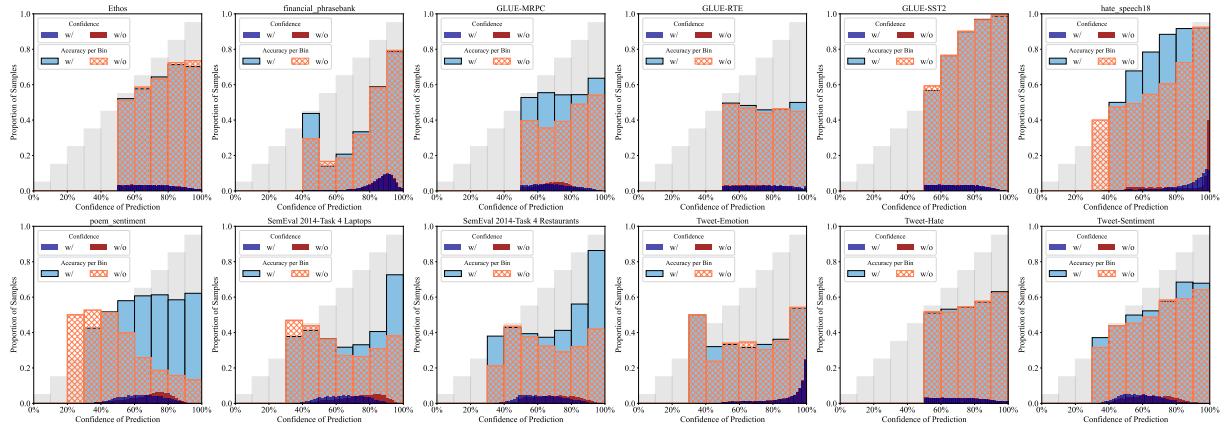


Figure 13: The reliability diagrams of GPT-J ( $k = 8$ ).

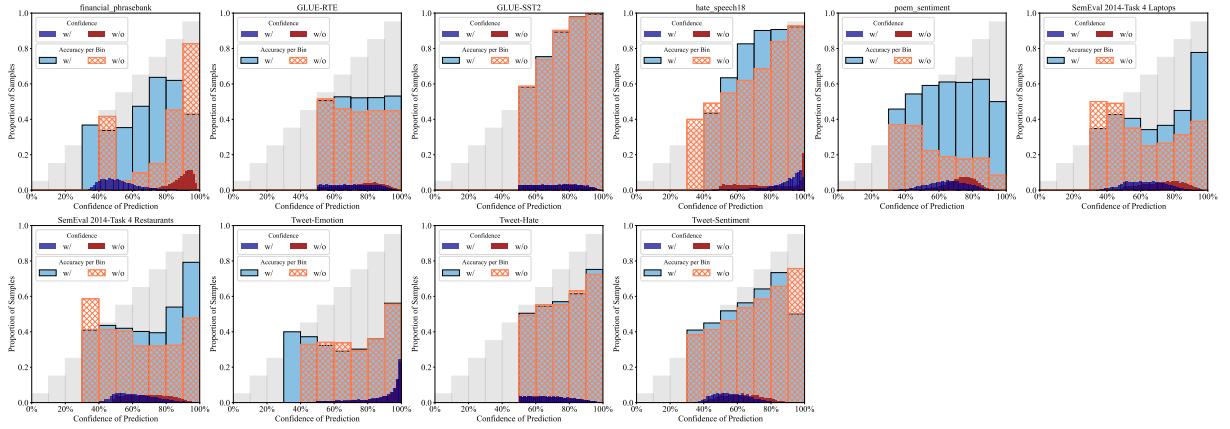


Figure 14: The reliability diagrams of GPT-J ( $k = 16$ ).

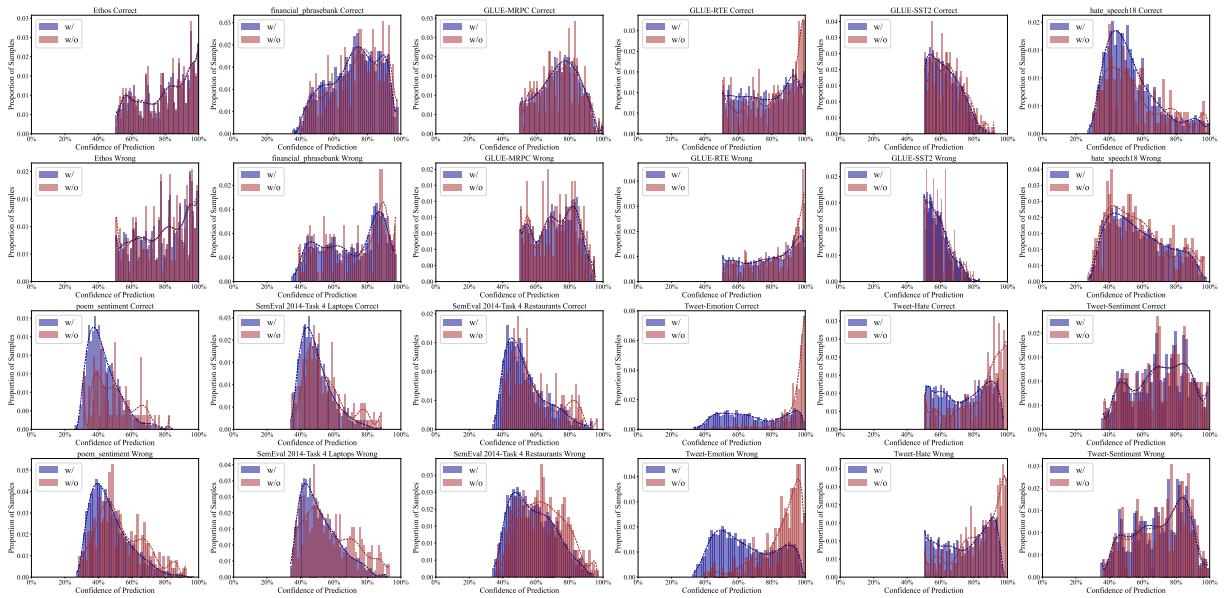


Figure 15: The confidence distribution of GPT-J ( $k = 0$ ).

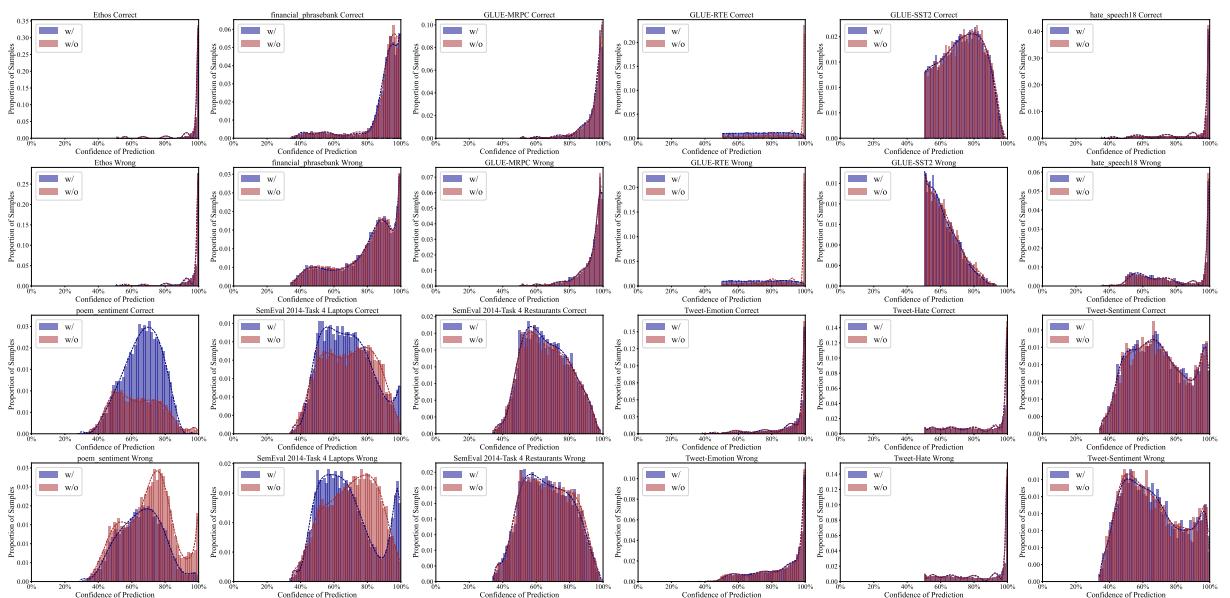


Figure 16: The confidence distribution of GPT-J ( $k = 1$ ).

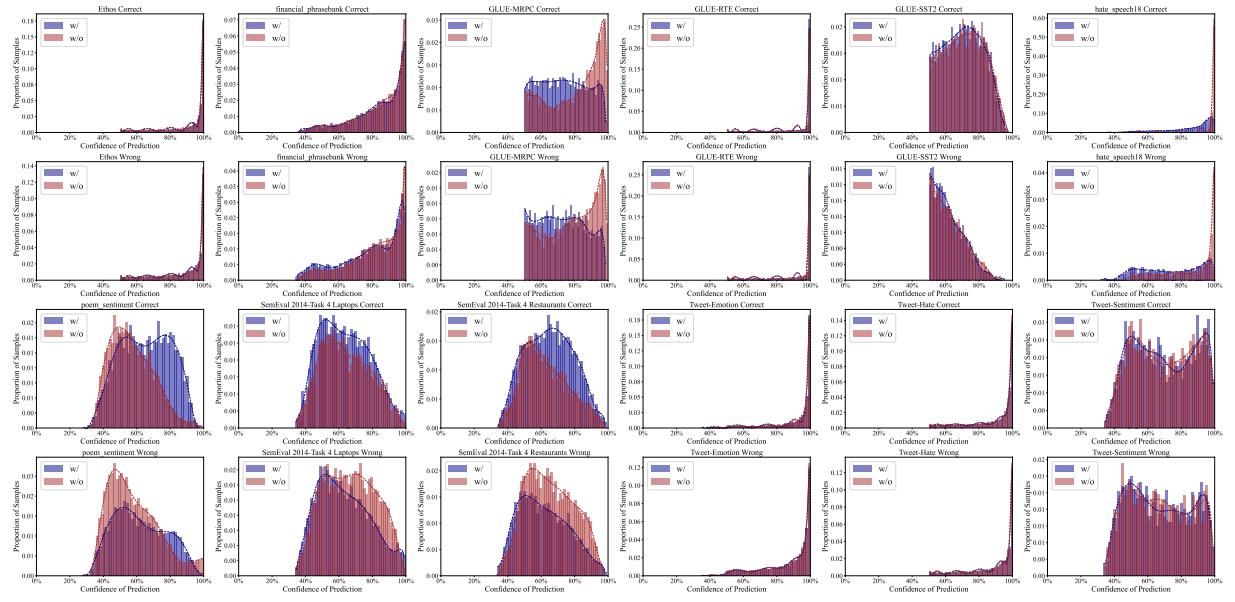


Figure 17: The confidence distribution of GPT-J ( $k = 2$ ).

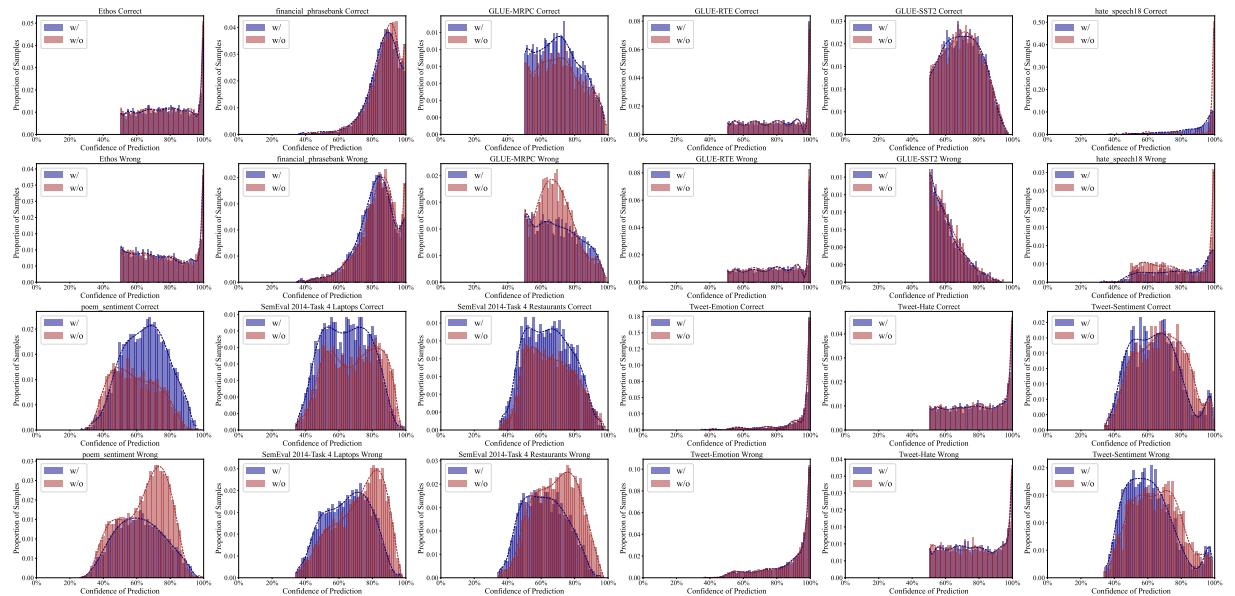


Figure 18: The confidence distribution of GPT-J ( $k = 4$ ).

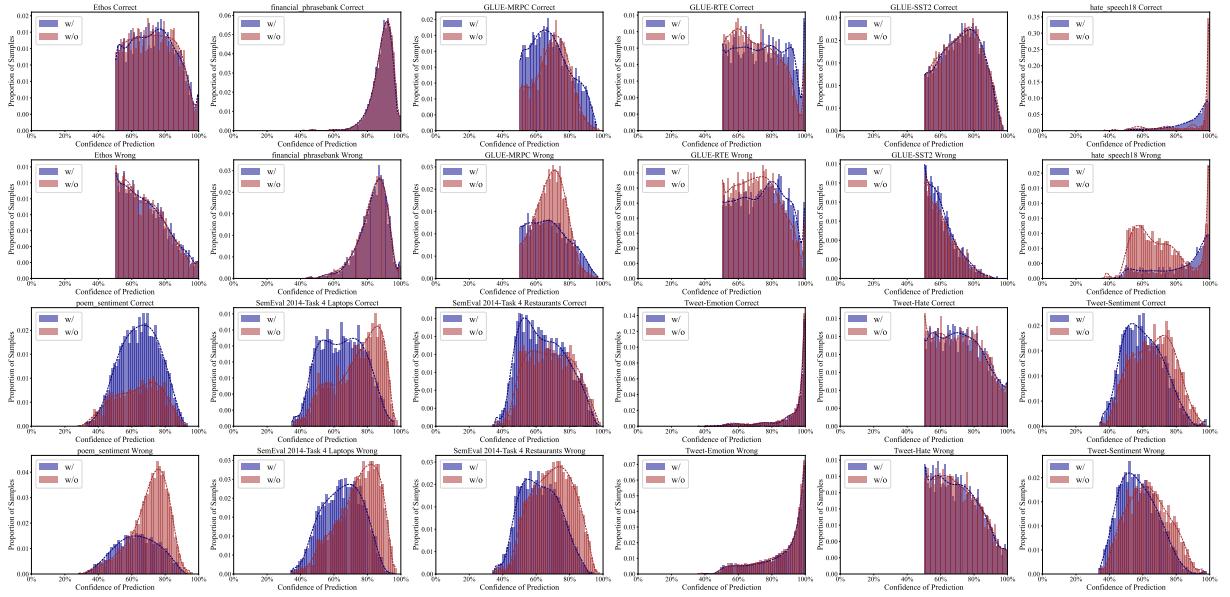


Figure 19: The confidence distribution of GPT-J ( $k = 8$ ).

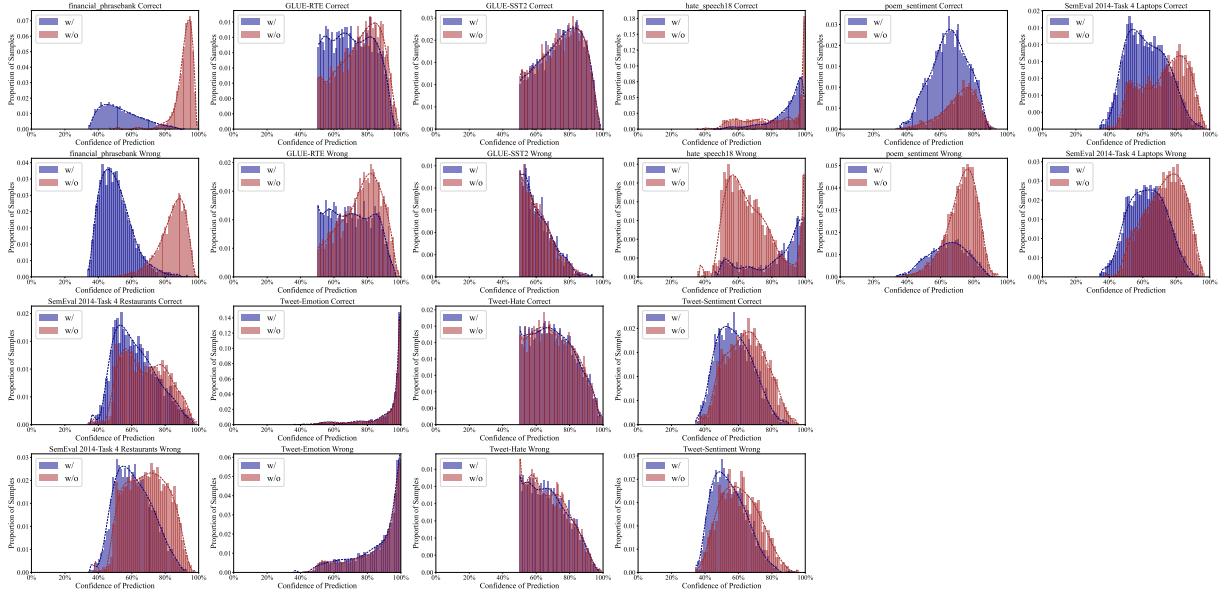


Figure 20: The confidence distribution of GPT-J ( $k = 16$ ).

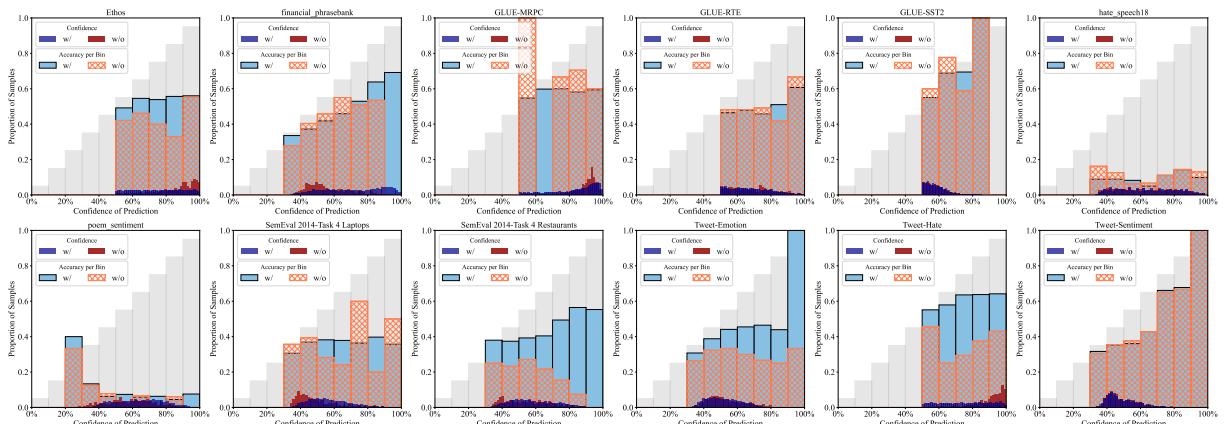


Figure 21: The reliability diagrams of GPT-2 ( $k = 0$ ).

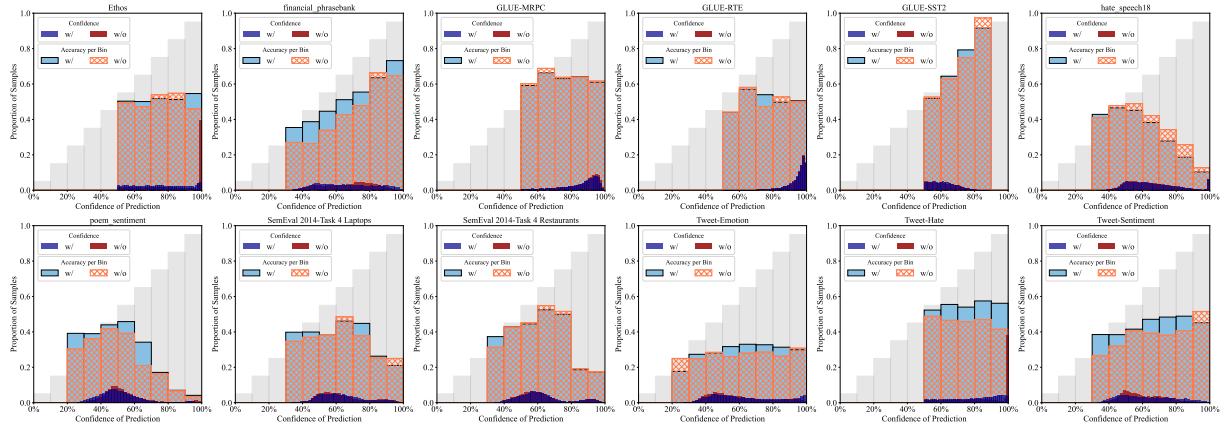


Figure 22: The reliability diagrams of GPT-2 ( $k = 1$ ).

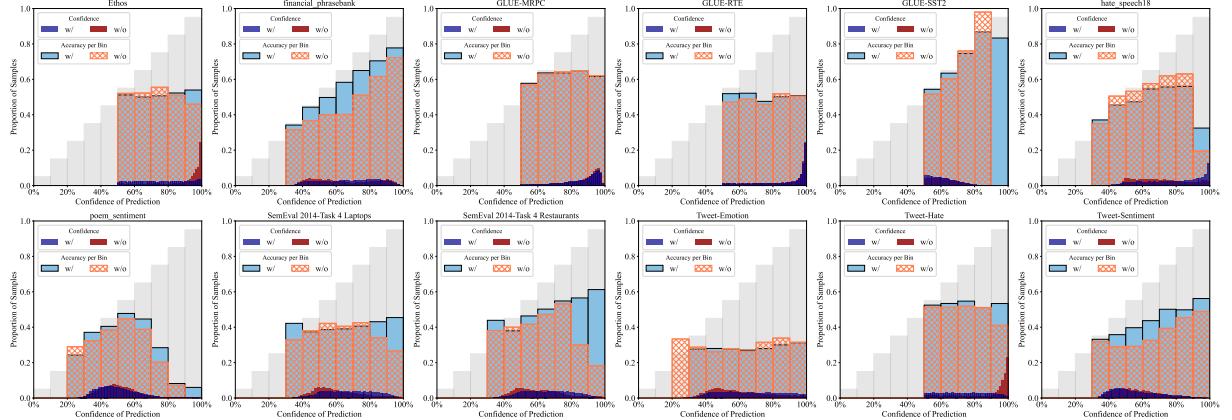


Figure 23: The reliability diagrams of GPT-2 ( $k = 2$ ).

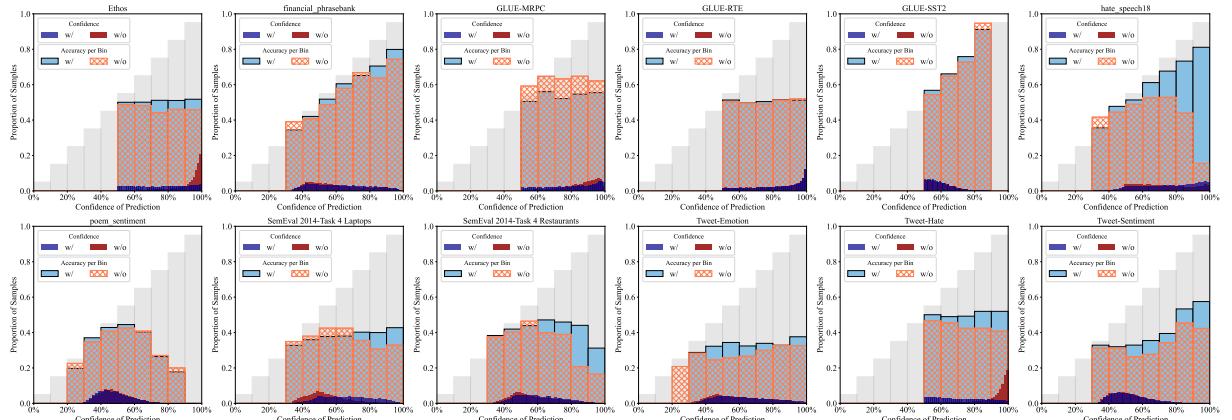


Figure 24: The reliability diagrams of GPT-2 ( $k = 4$ ).

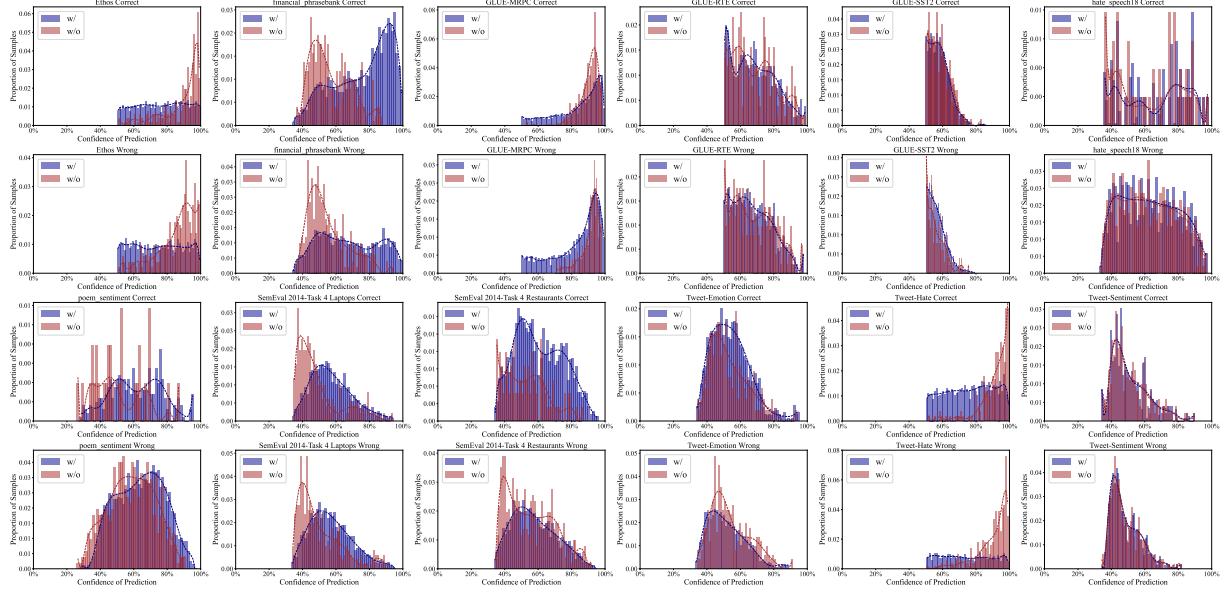


Figure 25: The confidence distribution of GPT-2 ( $k = 0$ ).

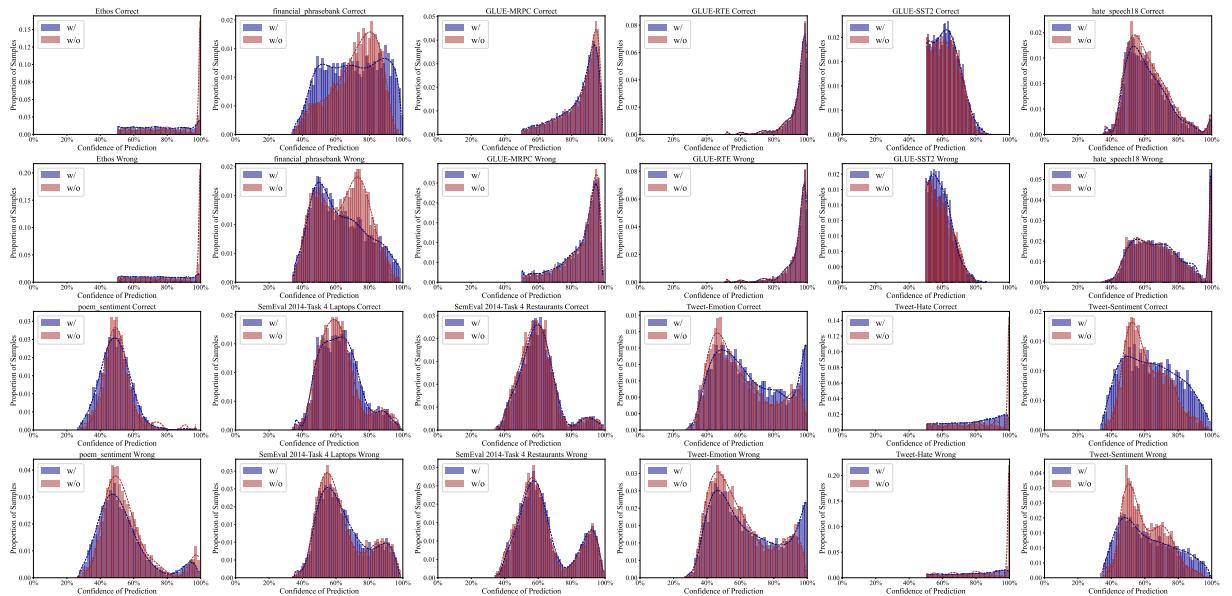


Figure 26: The confidence distribution of GPT-2 ( $k = 1$ ).

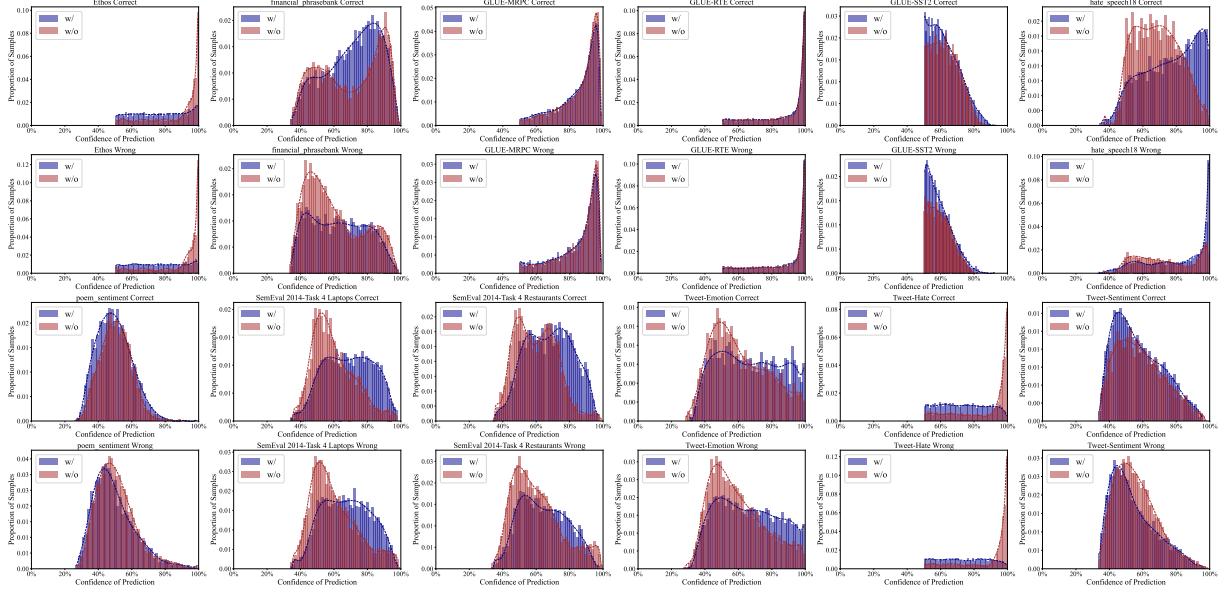


Figure 27: The confidence distribution of GPT-2 ( $k = 2$ ).

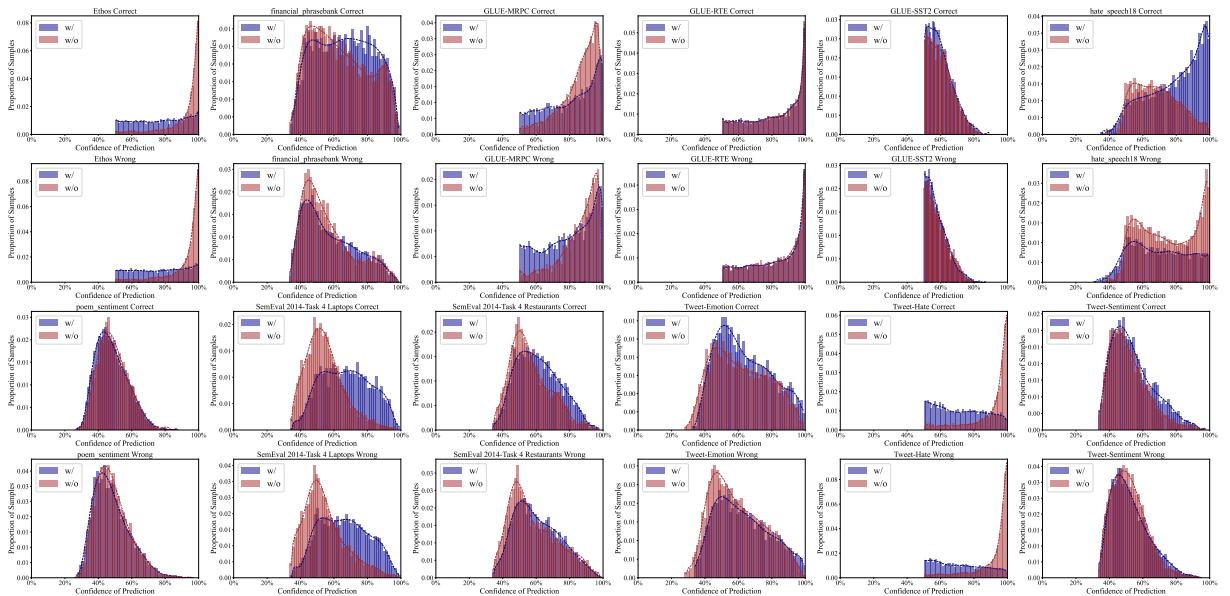


Figure 28: The confidence distribution of GPT-2 ( $k = 4$ ).