

On Distributional Robustness of In-Context Learning for Text Classification

Carolina Hatanpää¹ Noah A. Smith² Sachin Kumar³

Abstract

Pretrained language models (LMs) have been shown to be capable of few-shot learning via in-context demonstrations. In this work, we examine if in-context learning generalizes out-of-distribution. On several text classification tasks and considering different kinds of realistic demonstration vs. test distribution shifts—domain, adversarial, and dialectal—we find a significant drop in test accuracy compared to in-distribution performance. Moreover, we find that the accuracy is inversely proportional to the number of demonstrations. To address this issue, we explore different zero-shot prompting approaches which do not rely on demonstrations. With six open-source language model families, we show that zero-shot prompting techniques which verbalize the target distribution in the prompt are able to close the gap to in-domain few-shot classification performance.

1. Introduction

Language models trained on only raw text have been shown to learn in-context, that is, perform new tasks simply by conditioning on a handful of demonstrations (Brown et al., 2020). By drawing parallels to gradient-based optimization (Dai et al., 2023; Deutch et al., 2023), prior work has suggested that in-context learning (ICL) is prone to relying on superficial features in the training examples (Mueller et al., 2023b). Hence, it can be brittle to distribution shifts between the demonstrations and test examples with more demonstrations and larger models showing sharper drops in task performance (Tang et al., 2023).

With text classification as a case study, in this work, we investigate the utility of zero-shot inference methods in mitigating that issue, as they do not rely on any demonstrations. Recent studies have suggested that ICL serves as a way to

^{*}Equal contribution ¹Microsoft ²University of Washington ³The Ohio State University. Correspondence to: Sachin Kumar <kumar.1145@osu.edu>.

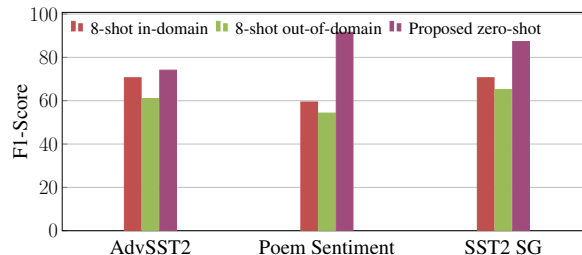


Figure 1. We observe considerable performance (F1) drops for sentiment classification when testing with three kinds of distribution shifts: adversarial, domain (poetry), and dialect (Singlish) when tested with Mistral 7B. The demonstrations are from SST2. Our proposed approach (§4) is able to close the gap.

implicitly prime the model with the domain, concepts, or topics and the format of the target task (Min et al., 2022; Wang et al., 2023). We build on the approach presented in Kumar et al. (2024) which showed that specifying these concepts explicitly in the prompt, without any demonstrations, can match or even surpass few-shot performance.

Evaluating on text classification tasks with three kinds of realistic distribution shifts (adversarial, domain, and dialectal) and with six language model families (model sizes ranging from 125M to 13B), we find that simply verbalizing the distributional properties of the test examples in the prompt, such as its domain or dialect, results in zero-shot approaches matching 8-shot in-domain performance. In a more realistic setting where the test distribution may not be known, we propose a simple but effective *mixture of prompts* approach that computes and aggregates model probabilities based on prompts containing verbalizations of different distributional shifts. Our experiments show that when the prompt mixture contains the true distribution information of the test example, this approach performs on par with the case where the properties are known in advance. In cases where the prompt mixture does not contain the test distribution description, the performance declines, however it still outperforms few-shot approaches.

2. Related Work

Since its introduction, various studies have attempted to analyze ICL’s underlying mechanisms (Xie et al., 2022;

On Distributional Robustness of In-Context Learning for Text Classification

	0 shot		1 shot		4 shot		8 shot		16 shot	
	ID	OOD	ID	OOD	ID	OOD	ID	OOD	ID	OOD
Poem	54.3 _{0.0}	54.3 _{0.0}	85.7 _{4.8}	85.7 _{13.3}	92.6 _{1.4}	80.0 _{9.9}	93.1 _{4.6}	52.6 _{3.4}	92.6 _{3.9}	88.0 _{2.8}
Finance	92.3 _{0.0}	92.3 _{0.0}	94.1 _{3.7}	95.1 _{0.6}	96.2 _{0.5}	90.4 _{5.0}	95.7 _{1.4}	74.3 _{32.5}	93.3 _{1.4}	93.3 _{5.1}

Table 1. Few-shot Macro-F1 (mean_{std}) with Pythia 6.9B for domain shifts with demonstrations from SST2. In-Distribution (ID) refers to cases where test cases and training data were from the same domain, and Out-of-Distribution (OOD) refers to cases where they were from different domains (see Appendix B for additional results).

	0 shot			1 shot			4 shot			8 shot		
	SAE	AAVE	SG	SAE	AAVE	SG	SAE	AAVE	SG	SAE	AAVE	SG
SST2	50.7 _{0.0}	52.9 _{0.0}	52.5 _{0.0}	88.9 _{3.2}	88.9 _{2.4}	86.4 _{2.6}	93.6 _{0.8}	90.8 _{2.1}	89.3 _{1.8}	93.6 _{1.3}	90.4 _{2.8}	88.9 _{3.3}
Emotion	42.0 _{0.0}	39.6 _{0.0}	38.7 _{0.0}	34.6 _{3.4}	33.5 _{4.2}	31.6 _{4.9}	31.8 _{7.5}	30.3 _{8.0}	29.6 _{8.8}	33.2 _{5.3}	32.3 _{5.8}	31.4 _{6.1}

Table 2. Few-shot Macro-F1 (mean_{std}) with Pythia 6.9B for dialect shifts from Standard American (SAE) to African American Vernacular (AAVE) and Singaporean English dialects (SG; see Appendix B for additional results).

Ahuja et al., 2023; Hahn & Goyal, 2023; Zhang et al., 2023; von Oswald et al., 2023; Wang et al., 2023) as well as explored its limits (Akyürek et al., 2022; Chan et al., 2022). Prior work has shown that ICL relies on shortcuts to make predictions (Tang et al., 2023). Saparov & He (2023) find that LMs tend to rely on spurious correlations for reasoning tasks that hinder robust generalization. However, they have largely focused on toy or synthetic tasks. In this work, we focus on realistic distribution shifts on popular text classification tasks such as sentiment, topic and emotion widely studied in the context of ICL. Most related to our work is Mueller et al. (2023a), which tests generalization for syntactic tasks.

While most prior work only points out the issue of robustness of ICL, we also explore potential solutions to tackle this issue using zero-shot approaches which in recent work have shown to match few-shot performance for text classification tasks Kumar et al. (2024). Relatedly, Drozdov et al. (2023), find that a series of chain-of-thought prompts can yield better robustness for semantic parsing.

3. ICL Is Not Distributionally Robust

In an in-context learning setup, given k demonstrations $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_k, y_k)\}$ and a test example \mathbf{x}_{test} , the label is predicted using a language model,

$$\max_{y \in \mathcal{Y}} p_{\text{LM}}(y \mid \mathbf{x}_1, y_1, \dots, \mathbf{x}_k, y_k, \mathbf{x}_{\text{test}}) \quad (1)$$

Here \mathcal{Y} is the set of all labels. In practice, the labels are verbalized in natural language (e.g., the words “negative” and “positive” for sentiment classification) and all the demonstrations are concatenated followed by the test example. Typically, all \mathbf{x}_i and $\mathbf{x}_{\text{test}} \sim \mathcal{D}$, an input distribution. Text classifiers *trained* with supervision have been shown to learn shortcuts in the data to make predictions which limit

their generalization. In this section, we examine this phenomenon for ICL by measuring if the classification accuracy is impacted when the test distribution diverges from the demonstrations, that is when $\mathbf{x}_{\text{test}} \sim \mathcal{D}'$ different from \mathcal{D} .

Experimental Setup and Results We evaluate on three tasks: sentiment (2-class), emotion (6-class), and topic (4-class) classification with demonstrations from SST2 (?movie reviews;][socher2013recursive, Emotions (Saravia et al., 2018), and AGNews (Zhang et al., 2015) respectively. We consider two kinds of distribution shifts in test examples: domain and dialectal. For domain shifts, we evaluate on two binary sentiment classification datasets: Poetry (Sheng & Uthus, 2020), and Financial News (Malo et al., 2014). We simulate the out-of-distribution (OOD) setup using demonstrations from the SST2 train set and compare with in-domain demonstrations. For dialectal shifts, we use multi-VALUE toolkit (Ziems et al., 2023) to translate SST2, Emotions, and AGNews to African American Vernacular English (AAVE) and Singaporean English (SG) using the original train set examples as demonstrations. We report macro F1-scores for all evaluation sets for $k \in \{0, 1, 4, 8\}$ where for each k we report mean and variance with 10 different sets of demonstrations sampled from the respective train sets. We evaluated on GPT2 (S, M, L, XL) (Radford et al., 2019), OPT (1.3B and 2.7B) (Zhang et al., 2022), Pythia (1.4B, 2.8B, and 6.9B) (Biderman et al., 2023), Mistral 7B (Jiang et al., 2023), Llama 1 (7B) (Touvron et al., 2023), and Llama 2 (7B and 13B) (Touvron et al., 2023). No finetuning was performed on these models.

We report results for Pythia 6.9B in Tables 1 and 2 (with results for other models following largely similar trends in Appendix B). We observe a significant decline in OOD accuracy for both domain and dialectal shifts with the former

being much larger. Furthermore, while more demonstrations increase in-domain performance, we observe an inverse trend with OOD setups.

4. Context Aware Zero-Shot Inference

As shown in Tables 1 and 2, the zero-shot setup observes no decline in performance across distributions since, by definition, it does not rely on any demonstrations. However, its overall accuracy is poor even in-domain compared to higher-shot setups. To close this gap, in this section, we describe *context-aware* zero-shot inference which uses the following inference objective: $\hat{y} = \max_{y \in \mathcal{Y}} p(y \mid \mathbf{x}_{\text{test}}, u, v, \dots)$. Here, u, v, \dots represent additional context that can aid the task (such as u could be domain of the input text, v could be the author demographic, and so on). Kumar et al. (2024), who proposed this approach, experimented with different instantiations of these factors showing that conditioning on the domain, subject, author, or reader of the input text can help improve classification performance to match up to 16-shot results while achieving low variance across prompt variations.

In this work, we explore the utility this setup where we consider the type of test distribution as additional context and verbalize it along with the label. We use a generative version of this setup, where $\hat{y} = \max_{y \in \mathcal{Y}} p(\mathbf{x}_{\text{test}} \mid y, u, v, \dots)$ to make predictions, which has been shown to perform better than the previously described discriminative setup. We verbalize the labels and the contextual factors (y, u, v, \dots) in textual form and refer to it as z (for example, “this is a positive^y movie review^u written in a Singaporean English dialect^v”). Furthermore, to reduce variance across prompt variations, we aggregate model probabilities across prompt paraphrases (Kumar et al., 2024) to obtain:

$$\hat{y}_{\text{GEN}} = \max_{y \in \mathcal{Y}} \sum_{z \in \mathcal{Z}(y, u, v, \dots)} p(\mathbf{x} \mid z) \quad (2)$$

\mathcal{Z} describes all potential ways to verbalize the label and the factors (in practice, we only use 10). However, the factors, i.e., the test distribution, might not be known. To overcome this issue, we propose to use a *mixture of prompts* which aggregates the probabilities in (2) over all *potential* factors,

$$\hat{y}_{\text{GEN}} = \max_{y \in \mathcal{Y}} \sum_{u, v, \dots} \sum_{z \in \mathcal{Z}(y, u, v, \dots)} p(\mathbf{x} \mid z) \quad (3)$$

The set of enumerable distributions can be infinitely large. In our experiments, we simulate a setting with a small set where we compare performance when the target distribution is either included or excluded from the set.

4.1. Experimental Setup

In addition to all the datasets described in §3, we evaluate on additional domains for sentiment classification, namely, CR (Customer Reviews; Hu & Liu, 2004), and MR (RottenTomatoes; Pang & Lee, 2005). We also experiment with an adversarial distribution shift by evaluating on adv-SST2 (Wang et al., 2022) (more details in Appendix A).

We consider three scenarios. The first is where we know the distribution of the test set (**known**). In this scenario, we make predictions using (2). In the second scenario, we do not know the distribution of the test set, but we assume that it exists in the set of potential distributions (**partly known**). In the third scenario, we have a completely unknown distribution that is not present in the set (**unknown**). For both of these scenarios, we use (3) for inference.

For each scenario, task and dataset, we hand write a label description z per label per distribution. We then generate 10 templatic paraphrases by querying ChatGPT (GPT-3.5) (Ouyang et al., 2022)¹ and manually verify the correctness of each paraphrase. For *mixture of prompt* experiments, for domain and adversarial shifts, we include the following distributions in the mixture: movie reviews, Amazon reviews, Yelp reviews, Poetry, Financial news, adversarial reviews. For dialectal shifts for sentiment, emotion and topic, we evaluate on AAVE and Singaporean English versions of the test sets and include AAVE, Singaporean, Indian, Irish, and Nigerian English in the mixture. Note that we are only including the names of these distributions in the prompts, and not any demonstrations. Finally, we simulate a case where we do not know the test distribution, by removing from the mixture, the prompts related to the test distribution. For each test set, we report mean and standard deviation of macro-F1.

In addition to the generative classification approach in (3), we also evaluate on the discriminative version defined by $p(z \mid \mathbf{x})$. We consider three versions of this approach which differ in how the context is provided (details in Appendix A)

For few-shot baselines, we consider two versions: one which uses (1) (ICL), and another which calibrates the label probability by dividing it with $p_{\text{LM}}(y \mid \text{NULL})$ (ICL-CC) (Zhao et al., 2021).

4.2. Results

We report representative results in Table 3 for Mistral 7B with the best results for each of the three zero-shot scenarios and the best performing few-shot baseline (with detailed

¹This process needs to be done only once for each task and, in practice, any paraphrasing model can be employed. We provide the list of all paraphrases we generated here: <https://pastebin.com/2gBYxJU>

	Few-shot(OOD)	Zero-shot(Known)	Zero-shot(Partly Known)	Zero-shot(Unknown)
Adv-SST2	62.7 _{5.8}	62.2 _{2.2}	73.9 _{0.6}	74.1 _{0.8}
MR	70.2 _{10.7}	90.2 _{0.2}	86.4 _{0.2}	81.3 _{0.2}
CR	71.5 _{8.1}	86.1 _{0.3}	87.8 _{0.4}	76.4 _{0.6}
Amazon	82.9 _{5.1}	92.3 _{0.3}	93.0 _{0.2}	92.3 _{0.1}
Yelp	82.5 _{6.8}	93.7 _{0.1}	92.7 _{0.1}	92.8 _{0.1}
Poem	59.4 _{6.9}	91.6 _{0.8}	89.6 _{0.8}	90.0 _{0.0}
Finance	90.6 _{0.0}	97.1 _{0.6}	91.3 _{0.9}	90.7 _{0.6}
SST2 AAVE	71.4 _{10.1}	88.8 _{0.3}	89.1 _{0.1}	89.3 _{0.2}
SST2 SG	69.3 _{9.1}	87.3 _{0.2}	88.6 _{0.2}	88.6 _{0.3}
Emotion AAVE	40.9 _{2.6}	51.4 _{0.2}	49.6 _{0.2}	48.0 _{0.1}
Emotions SG	40.8 _{2.5}	48.0 _{0.2}	46.5 _{0.1}	49.2 _{0.1}
AGNews AAVE	82.6 _{0.0}	48.1 _{0.1}	48.1 _{0.0}	47.5 _{0.0}
AGNews SG	89.9 _{8.2}	48.3 _{0.1}	47.6 _{0.0}	47.4 _{0.0}

Table 3. Macro-F1 with Mistral 7B. We report mean_{std} over 10 runs for zero-shot approaches over 5 seeds for the few-shot ones. For each result the best approach within that category is shown; Detailed results in Appendix B.

breakdown for all models with largely similar trends in Appendix B).

We find that if the test distribution is known, the zero-shot approaches achieve the best performance overall surpassing the few-shot out-of-distribution accuracy by a large margin and even approaching the in-distribution test performance (see Figure 1). This trend holds for the majority of the 12 models we tested (in 8-10 cases depending on the dataset).

In a more realistic scenario, if the test distribution is unknown, the mixture of prompts approach described in (3) (partially unknown) still performs better than the few-shot baselines. However, the improvements are not as large compared to the first scenario as explicitly priming the model with a different distribution than the target might hurt. When the target distribution is not included in the prompt mixture (unknown), we find that this approach still performs better than baseline, however, the accuracy gain is lower than previous two scenarios. While “unknown” depicts the worst case scenario, for most practical purposes, given a large enough set of distributions, the target distribution is highly likely to be included in the set.

Finally, the generative zero-shot classification objective is the winning approach in most settings.² However, the discriminative versions are competitive and even out-perform the generative approach in a few cases (see details in Appendix A).

²The AGNews dataset does not reflect these results. This is likely because our dialectal data creation method does not modify topical words in the test examples.

5. Conclusion

In this work, we investigate the out-of-distribution robustness of in-context learning. On several text classification tasks we find that a distribution mismatch between the test instances and the demonstrations can lead to a considerable drop in accuracy compared to in-distribution accuracy. We propose to address this issue using zero-shot prompting. Our experiments show that simply describing the distribution in the prompt text can considerably close the performance gap.

6. Limitations

The datasets in our study are exclusively in English, so the generalizability of our paper’s findings to other languages may be limited. In addition, our dialectal distribution datasets are synthetically generated, so they may not perfectly reflect the dialects they represent. We also make simplifying assumptions about label independence which may not be true in practice. Due to computing limitations, we were unable to run experiments on larger models, which may not indicate the same trends.

References

- Ahuja, K., Panwar, M., and Goyal, N. In-context learning through the bayesian prism, 2023.
- Akyürek, E., Schuurmans, D., Andreas, J., Ma, T., and Zhou, D. What learning algorithm is in-context learning? investigations with linear models. *arXiv preprint arXiv:2211.15661*, 2022.
- Biderman, S., Schoelkopf, H., Anthony, Q. G., Bradley, H., O’Brien, K., Hallahan, E., Khan, M. A., Purohit, S.,

- Prashanth, U. S., Raff, E., et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- Chan, S. C., Dasgupta, I., Kim, J., Kumaran, D., Lampinen, A. K., and Hill, F. Transformers generalize differently from information stored in context vs in weights. *arXiv preprint arXiv:2210.05675*, 2022.
- Dai, D., Sun, Y., Dong, L., Hao, Y., Ma, S., Sui, Z., and Wei, F. Why can GPT learn in-context? language models secretly perform gradient descent as meta-optimizers. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 4005–4019, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.247. URL <https://aclanthology.org/2023.findings-acl.247>.
- Deutch, G., Magar, N., Bar Natan, T., and Dar, G. In-context learning and gradient descent revisited. *arXiv e-prints*, pp. arXiv–2311, 2023.
- Drozdo, A., Schärli, N., Akyürek, E., Scales, N., Song, X., Chen, X., Bousquet, O., and Zhou, D. Compositional semantic parsing with large language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=gJW8hSGBys8>.
- Hahn, M. and Goyal, N. A theory of emergent in-context learning as implicit structure induction, 2023.
- Hu, M. and Liu, B. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 168–177, 2004.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mistral 7b, 2023.
- Kumar, S., Park, C. Y., and Tsvelkov, Y. Gen-z: Generative zero-shot text classification with contextualized label descriptions. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=rkplyfqr0>.
- Malo, P., Sinha, A., Korhonen, P., Wallenius, J., and Takala, P. Good debt or bad debt: Detecting semantic orientations in economic texts. *Journal of the Association for Information Science and Technology*, 65(4):782–796, 2014.
- Min, S., Lyu, X., Holtzman, A., Artetxe, M., Lewis, M., Hajishirzi, H., and Zettlemoyer, L. Rethinking the role of demonstrations: What makes in-context learning work? In Goldberg, Y., Kozareva, Z., and Zhang, Y. (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 11048–11064, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.759. URL <https://aclanthology.org/2022.emnlp-main.759>.
- Mueller, A., Narang, K., Mathias, L., Wang, Q., and Firooz, H. Meta-training with demonstration retrieval for efficient few-shot learning. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 6049–6064, Toronto, Canada, July 2023a. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.376. URL <https://aclanthology.org/2023.findings-acl.376>.
- Mueller, A., Webson, A., Petty, J., and Linzen, T. In-context learning generalizes, but not always robustly: The case of syntax. *ArXiv*, abs/2311.07811, 2023b. URL <https://api.semanticscholar.org/CorpusID:265158068>.
- Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Pang, B. and Lee, L. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. *arXiv preprint cs/0506075*, 2005.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Saparov, A. and He, H. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought, 2023.
- Saravia, E., Liu, H.-C. T., Huang, Y.-H., Wu, J., and Chen, Y.-S. Carer: Contextualized affect representations for emotion recognition. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pp. 3687–3697, 2018.

- Sheng, E. and Uthus, D. Investigating societal biases in a poetry composition system. *arXiv preprint arXiv:2011.02686*, 2020.
- Tang, R., Kong, D., Huang, L., and Xue, H. Large language models can be lazy learners: Analyze shortcuts in in-context learning. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 4645–4657, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.284. URL <https://aclanthology.org/2023.findings-acl.284>.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- von Oswald, J., Niklasson, E., Randazzo, E., Sacramento, J., Mordvintsev, A., Zhmoginov, A., and Vladymyrov, M. Transformers learn in-context by gradient descent, 2023.
- Wang, B., Xu, C., Wang, S., Gan, Z., Cheng, Y., Gao, J., Awadallah, A. H., and Li, B. Adversarial glue: A multi-task benchmark for robustness evaluation of language models, 2022.
- Wang, X., Zhu, W., Saxon, M., Steyvers, M., and Wang, W. Y. Large language models are implicitly topic models: Explaining and finding good demonstrations for in-context learning. In *Workshop on Efficient Systems for Foundation Models@ ICML2023*, 2023.
- Xie, S. M., Raghunathan, A., Liang, P., and Ma, T. An explanation of in-context learning as implicit bayesian inference. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=RdJVFCHjUMI>.
- Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X. V., et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- Zhang, X., Zhao, J., and LeCun, Y. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28, 2015.
- Zhang, Y., Zhang, F., Yang, Z., and Wang, Z. What and how does in-context learning learn? bayesian model averaging, parameterization, and generalization, 2023.
- Zhao, Z., Wallace, E., Feng, S., Klein, D., and Singh, S. Calibrate before use: Improving few-shot performance of language models. In *International conference on machine learning*, pp. 12697–12706. PMLR, 2021.
- Ziems, C., Held, W., Yang, J., Dhamala, J., Gupta, R., and Yang, D. Multi-VALUE: A framework for cross-dialectal English NLP. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 744–768, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.44. URL <https://aclanthology.org/2023.acl-long.44>.

A. Additional Experimental Details

Discriminative Baselines:

We consider three variations of each discriminative zero-shot baseline (see Table 18): (1) Direct context specifies the contextual variable using a simple format, (2) Direct instruct specifies the contextual variable and also provides the labels the model is expected to predict, and (3) Direct simple does not provide the contextual variable. Table 4 summarizes all of the datasets we use. Table 17 summarizes the handwritten label description templates we use and paraphrase to construct the full set of label descriptions.

B. Detailed Results

Detailed results for the following models are provided in this Appendix: GPT2 Results in Table 5, GPT2-Medium Results in Table 6, GPT2-Large Results in Table 7, GPT2-XL Results in Table 8, OPT 1.3B Results in Table 9, Pythia 1.4B Results in Table 10, OPT 2.7B Results in Table 11, Pythia 2.8B Results in Table 12, Pythia 6.9B Results in Table 13, Mistral 7B Results in Table 14, Llama 7B Results in Table 15, and Llama-2 7B Results in Table 16.

Dataset	Task	Domain	# of Classes	Inclusion in Minimal Set
SST2	Sentiment Classification	Movies	2	Yes
ADV SST2	Sentiment Classification	Movies, adversarial	2	Yes
MR	Sentiment Classification	Movies, Rotten Tomatoes	2	No
CR	Sentiment Classification	Customer Reviews	2	Yes
Amazon	Sentiment Classification	Customer Reviews, Amazon	2	No
Yelp	Sentiment Classification	Yelp reviews	2	Yes
Poem	Sentiment Classification	Poetry	4, reduced to 2	Yes
Finance	Sentiment Classification	Economic News	3, reduced to 2	Yes
SST2 AAVE	Sentiment Classification	Movies, African American Vernacular English (AAVE), synthetically created	2	Yes
SST2 SG	Sentiment Classification	movies, Singaporean English (Singlish), synthetically created	2	Yes
Emotion AAVE	Emotion Classification	Tweets, African American Vernacular English (AAVE), synthetically created	6	Yes
Emotion SG	Emotion Classification	Tweets, Singaporean English (Singlish), synthetically created	6	Yes
AG News AAVE	Topic Classification	News, African American Vernacular English (AAVE), synthetically created	4	Yes
AG News SG	Topic Classification	News, Singaporean English (Singlish), synthetically created	4	Yes

Table 4. Datasets used for the experiments

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	83.1 _{0.3}	72.8 _{0.5}	58.4 _{0.3}	74.5 _{0.6}	68.3 _{0.3}	73.9 _{0.2}	77.6 _{0.2}	73.6 _{0.5}	64.5 _{0.4}	72.7 _{0.4}	61.0 _{0.4}	70.2 _{0.2}	66.5 _{0.0}	57.5 _{3.3}	58.6 _{7.8}	55.9 _{6.1}
SST2 ADV	42.1 _{1.7}	46.6 _{1.1}	43.3 _{2.1}	44.7 _{1.4}	48.1 _{0.7}	50.9 _{0.8}	43.1 _{0.4}	43.6 _{0.6}	48.2 _{0.6}	49.5 _{1.4}	44.6 _{0.6}	42.7 _{0.8}	44.6 _{0.0}	47.4 _{3.3}	49.2 _{1.7}	46.9 _{2.3}
MR	76.9 _{0.3}	66.2 _{0.6}	53.8 _{0.5}	64.8 _{0.6}	66.3 _{0.5}	70.0 _{0.2}	72.6 _{0.2}	70.1 _{0.4}	63.1 _{0.3}	66.9 _{0.2}	58.8 _{0.4}	65.8 _{0.4}	62.4 _{0.0}	55.5 _{4.0}	53.1 _{2.7}	55.7 _{6.8}
CR	73.4 _{0.3}	69.2 _{0.7}	40.0 _{0.4}	55.1 _{0.9}	63.6 _{0.7}	61.9 _{0.4}	45.1 _{0.2}	69.4 _{0.7}	63.2 _{0.8}	63.0 _{0.7}	45.0 _{0.4}	71.1 _{0.6}	73.4 _{0.0}	54.0 _{6.0}	40.4 _{4.3}	45.1 _{11.8}
Amazon	78.1 _{0.4}	77.0 _{0.6}	62.4 _{0.6}	73.4 _{0.6}	76.5 _{0.7}	74.5 _{0.3}	84.1 _{0.1}	76.2 _{0.3}	76.6 _{0.4}	75.2 _{0.3}	84.0 _{0.1}	77.3 _{0.3}	67.1 _{0.0}	56.5 _{2.3}	57.9 _{3.4}	61.4 _{6.6}
Yelp	80.4 _{0.3}	82.5 _{0.5}	62.8 _{0.3}	82.2 _{0.4}	77.7 _{0.2}	77.6 _{0.1}	82.5 _{0.2}	76.9 _{0.2}	78.2 _{0.3}	76.8 _{0.2}	82.6 _{0.2}	77.7 _{0.2}	75.1 _{0.0}	54.3 _{4.4}	57.2 _{3.0}	62.8 _{9.8}
Poem	85.6 _{2.3}	52.8 _{3.7}	48.0 _{0.0}	56.2 _{4.3}	78.2 _{2.2}	74.6 _{1.0}	58.0 _{0.0}	68.4 _{0.8}	79.4 _{3.0}	75.8 _{2.2}	58.0 _{0.0}	68.0 _{0.0}	77.1 _{0.0}	51.4 _{6.0}	57.7 _{4.2}	54.3 _{0.0}
Finance	79.8 _{1.0}	75.1 _{1.6}	47.5 _{2.3}	72.0 _{1.2}	71.4 _{0.7}	59.5 _{0.9}	40.9 _{0.5}	65.3 _{1.2}	70.8 _{0.6}	58.9 _{0.7}	41.3 _{0.4}	65.0 _{0.4}	50.2 _{0.0}	29.6 _{7.0}	47.7 _{24.8}	90.6 _{0.0}
SST2 AAVE	80.3 _{0.3}	54.5 _{0.4}	60.5 _{0.6}	65.9 _{0.6}	80.4 _{0.3}	50.9 _{0.1}	50.9 _{0.1}	51.0 _{0.1}	80.1 _{0.2}	58.4 _{0.6}	59.4 _{1.0}	67.5 _{0.3}	64.0 _{0.0}	55.1 _{4.3}	52.2 _{0.9}	49.1 _{0.0}
SST2 SG	78.2 _{0.4}	59.6 _{0.6}	51.5 _{0.2}	67.2 _{0.6}	77.9 _{0.2}	50.9 _{0.1}	50.9 _{0.1}	50.9 _{0.1}	78.1 _{0.2}	56.9 _{0.3}	62.8 _{0.3}	65.7 _{0.4}	64.1 _{0.0}	56.1 _{3.5}	51.6 _{2.0}	53.9 _{6.0}

 Table 5. Macro-F1 with GPT2. We report mean_{std} over 10 runs for zero-shot. For few-shot, we report *average_{std}* over 5 seeds.

On Distributional Robustness of In-Context Learning for Text Classification

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	88.6 _{3.2}	76.6 _{0.4}	81.5 _{0.5}	73.1 _{0.5}	80.1 _{0.3}	77.3 _{0.1}	80.0 _{0.4}	77.6 _{0.3}	76.6 _{0.5}	83.3 _{0.2}	75.3 _{0.2}	81.3 _{0.3}	77.0 _{0.0}	60.1 _{9.9}	58.3 _{13.3}	63.2 _{12.3}
SST2 ADV	44.2 _{1.5}	44.0 _{1.0}	48.9 _{1.4}	42.2 _{2.1}	51.5 _{1.2}	47.0 _{0.9}	51.3 _{0.8}	44.8 _{0.8}	50.6 _{0.8}	45.7 _{0.6}	50.1 _{0.5}	44.2 _{0.8}	43.9 _{0.0}	45.4 _{3.4}	48.0 _{2.1}	48.8 _{3.2}
MR	80.1 _{0.30}	70.0 _{1.0}	79.2 _{0.6}	66.5 _{0.8}	75.0 _{0.3}	73.8 _{0.3}	77.1 _{0.2}	77.8 _{0.2}	71.4 _{0.3}	80.1 _{0.3}	73.7 _{0.3}	78.4 _{0.2}	73.9 _{0.0}	58.7 _{8.7}	57.6 _{11.9}	63.2 _{10.4}
CR	84.0 _{0.3}	85.1 _{1.0}	51.4 _{1.0}	73.2 _{1.3}	78.3 _{0.6}	61.6 _{0.9}	45.2 _{0.2}	63.9 _{0.4}	78.8 _{0.6}	61.2 _{1.0}	45.0 _{0.3}	64.0 _{0.8}	78.2 _{0.0}	62.3 _{11.4}	52.9 _{17.3}	74.3 _{6.7}
Amazon	62.0 _{0.3}	86.7 _{0.4}	75.3 _{0.5}	83.2 _{0.4}	75.1 _{0.3}	80.5 _{0.3}	80.9 _{0.1}	82.1 _{0.3}	76.3 _{0.4}	80.8 _{0.2}	80.7 _{0.1}	82.4 _{0.3}	75.7 _{0.0}	61.1 _{9.1}	65.7 _{8.9}	66.8 _{16.1}
Yelp	85.8 _{0.2}	90.3 _{0.1}	73.2 _{0.4}	89.1 _{0.2}	80.7 _{0.2}	77.6 _{0.2}	81.3 _{0.1}	77.2 _{0.2}	80.3 _{0.3}	77.5 _{0.2}	81.2 _{0.1}	77.1 _{0.2}	87.9 _{0.0}	65.4 _{11.0}	70.6 _{10.9}	67.2 _{17.8}
Poem	84.8 _{1.7}	63.0 _{2.2}	48.2 _{0.6}	77.6 _{2.5}	71.4 _{2.1}	77.4 _{1.0}	80.0 _{0.0}	76.0 _{0.9}	72.0 _{1.9}	75.8 _{0.6}	80.0 _{0.0}	76.0 _{1.3}	77.1 _{0.0}	59.4 _{7.1}	54.3 _{6.5}	54.3 _{0.0}
Finance	90.0 _{0.5}	90.7 _{1.4}	59.9 _{2.4}	86.1 _{1.0}	82.7 _{0.4}	77.5 _{0.6}	84.2 _{0.3}	71.0 _{0.9}	82.4 _{0.9}	76.8 _{0.9}	85.0 _{0.9}	73.8 _{0.8}	86.5 _{0.0}	74.6 _{32.6}	68.2 _{31.5}	90.6 _{0.0}
SST2 AAVE	83.4 _{0.5}	71.6 _{0.6}	72.4 _{0.4}	75.7 _{0.8}	83.8 _{0.2}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	82.4 _{0.2}	75.2 _{0.3}	53.7 _{0.2}	73.5 _{0.2}	75.1 _{0.0}	60.1 _{9.0}	57.5 _{12.9}	63.4 _{10.9}
SST2 SG	79.0 _{0.3}	64.1 _{0.7}	59.2 _{0.4}	73.7 _{0.4}	81.0 _{0.2}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	81.4 _{0.1}	70.6 _{0.4}	59.0 _{0.4}	69.9 _{0.4}	74.8 _{0.0}	61.0 _{8.9}	56.4 _{12.0}	62.0 _{9.6}

Table 6. Macro-F1 with GPT2-Medium. We report $average_{std}$ over 10 runs for zero-shot. For few-shot, we report $average_{std}$ over 5 seeds.

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	88.6 _{0.2}	84.3 _{0.5}	79.9 _{0.3}	83.5 _{0.4}	81.6 _{0.3}	83.8 _{0.4}	76.6 _{0.1}	83.1 _{0.2}	75.2 _{0.3}	83.6 _{0.1}	75.0 _{0.3}	83.6 _{0.2}	69.8 _{0.0}	74.1 _{11.5}	73.7 _{10.1}	81.0 _{9.5}
SST2 ADV	47.6 _{0.9}	51.1 _{1.2}	58.9 _{1.1}	51.2 _{1.8}	52.3 _{1.0}	47.8 _{1.0}	51.0 _{0.4}	51.1 _{1.0}	52.9 _{0.9}	47.5 _{0.8}	51.6 _{0.9}	53.4 _{0.7}	43.2 _{0.0}	41.5 _{5.4}	52.6 _{2.5}	49.0 _{3.3}
MR	82.9 _{0.2}	72.4 _{0.6}	78.0 _{0.5}	72.2 _{0.5}	76.5 _{0.3}	78.2 _{0.3}	75.6 _{0.1}	78.6 _{0.2}	71.8 _{0.5}	80.6 _{0.2}	71.5 _{0.5}	79.2 _{0.3}	76.2 _{0.0}	72.5 _{11.1}	67.5 _{9.2}	80.3 _{10.7}
CR	82.6 _{0.4}	83.4 _{0.7}	74.1 _{1.1}	80.4 _{0.5}	73.9 _{0.8}	77.9 _{0.5}	47.7 _{0.2}	79.0 _{0.8}	72.9 _{0.7}	72.8 _{0.4}	47.6 _{0.2}	76.2 _{0.8}	84.8 _{0.0}	76.3 _{6.4}	71.7 _{4.6}	78.3 _{6.7}
Amazon	84.6 _{0.4}	86.5 _{0.4}	83.1 _{0.5}	80.0 _{0.5}	85.8 _{0.4}	86.9 _{0.2}	74.7 _{0.1}	87.9 _{0.3}	85.8 _{0.4}	86.6 _{0.2}	74.7 _{0.1}	87.3 _{0.3}	90.6 _{0.0}	80.9 _{13.5}	79.8 _{3.8}	81.5 _{5.3}
Yelp	87.1 _{0.2}	90.5 _{0.3}	80.2 _{0.4}	89.1 _{0.4}	85.3 _{0.3}	75.9 _{0.2}	70.9 _{0.2}	76.0 _{0.3}	85.0 _{0.3}	74.6 _{0.2}	70.9 _{0.1}	74.9 _{0.3}	94.7 _{0.0}	80.3 _{15.2}	76.3 _{3.2}	74.6 _{7.1}
Poem	86.6 _{1.6}	80.6 _{2.5}	46.4 _{0.8}	84.6 _{3.3}	70.6 _{2.3}	84.0 _{1.6}	52.0 _{0.9}	77.8 _{1.1}	73.2 _{1.7}	84.6 _{1.3}	51.6 _{1.3}	77.8 _{1.1}	80.0 _{0.0}	56.0 _{5.9}	45.7 _{4.8}	52.6 _{3.4}
Finance	85.7 _{0.9}	92.6 _{0.9}	37.4 _{0.9}	83.3 _{1.1}	73.2 _{0.6}	85.6 _{0.4}	47.7 _{0.6}	82.5 _{0.8}	73.1 _{0.4}	87.8 _{0.6}	47.6 _{0.6}	86.2 _{0.4}	92.3 _{0.0}	70.4 _{31.7}	61.1 _{30.1}	74.3 _{32.5}
SST2 AAVE	83.8 _{0.3}	52.5 _{0.2}	50.7 _{0.2}	82.3 _{0.8}	83.5 _{0.1}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	84.3 _{0.3}	78.1 _{0.3}	57.1 _{0.3}	80.9 _{0.4}	76.2 _{0.0}	72.6 _{10.7}	67.7 _{10.0}	79.8 _{11.7}
SST2 SG	81.7 _{0.2}	55.9 _{0.4}	51.3 _{0.2}	78.4 _{0.4}	84.1 _{0.2}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	83.9 _{0.1}	76.5 _{0.6}	53.2 _{0.1}	80.6 _{0.4}	73.3 _{0.0}	73.3 _{11.1}	66.3 _{10.3}	79.2 _{11.9}

Table 7. Macro-F1 with GPT2-Large. We report $average_{std}$ over 10 runs for zero-shot. For few-shot, we report $average_{std}$ over 5 seeds.

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	90.3 _{0.2}	88.8 _{0.3}	70.6 _{0.2}	87.2 _{0.3}	83.4 _{0.3}	85.2 _{0.2}	78.6 _{0.2}	85.3 _{0.2}	77.8 _{0.4}	79.2 _{0.3}	73.5 _{0.4}	77.6 _{0.5}	64.8 _{0.0}	78.3 _{1.6}	79.3 _{10.8}	72.5 _{11.1}
SST2 ADV	54.7 _{0.8}	43.1 _{1.8}	49.1 _{1.0}	47.2 _{1.7}	53.6 _{1.0}	58.5 _{0.8}	53.1 _{0.5}	61.1 _{1.1}	54.3 _{1.1}	57.4 _{0.8}	53.0 _{0.7}	61.2 _{0.6}	48.7 _{0.0}	50.4 _{2.0}	51.8 _{6.9}	55.8 _{7.0}
MR	84.2 _{0.1}	87.3 _{0.3}	64.2 _{0.4}	84.6 _{0.6}	77.5 _{0.3}	82.8 _{0.3}	75.4 _{0.2}	83.1 _{0.2}	73.0 _{0.3}	75.6 _{0.3}	69.7 _{0.4}	76.6 _{0.4}	68.2 _{0.0}	74.9 _{2.0}	77.4 _{8.7}	71.1 _{9.9}
CR	86.2 _{0.3}	90.5 _{0.5}	41.2 _{0.9}	85.4 _{0.5}	84.6 _{0.6}	78.0 _{0.3}	70.6 _{0.9}	77.8 _{0.3}	84.7 _{0.6}	77.2 _{0.4}	70.3 _{0.6}	77.6 _{0.5}	81.7 _{0.0}	79.1 _{76.1}	80.1 _{5.6}	76.2 _{6.7}
Amazon	90.3 _{0.3}	90.9 _{0.4}	61.5 _{0.4}	79.5 _{0.5}	87.6 _{0.3}	85.2 _{0.3}	86.6 _{0.2}	85.5 _{0.2}	87.5 _{0.3}	85.0 _{0.3}	86.5 _{0.1}	85.4 _{0.3}	83.3 _{0.0}	86.3 _{9.6}	83.8 _{4.0}	80.3 _{8.0}
Yelp	88.3 _{0.1}	91.9 _{0.2}	61.7 _{0.3}	91.5 _{0.2}	86.9 _{0.1}	76.4 _{0.3}	87.9 _{0.1}	75.1 _{0.3}	86.7 _{0.2}	76.5 _{0.3}	87.9 _{0.2}	75.8 _{0.3}	85.1 _{0.0}	85.4 _{12.6}	83.2 _{7.1}	80.2 _{7.7}
Poem	82.2 _{2.2}	78.0 _{2.1}	48.4 _{0.8}	76.6 _{3.7}	80.6 _{1.3}	85.4 _{1.9}	70.6 _{1.3}	79.6 _{1.3}	82.6 _{1.9}	84.0 _{0.9}	70.2 _{0.3}	79.0 _{1.1}	57.1 _{0.0}	70.9 _{7.1}	68.0 _{9.3}	52.6 _{3.4}
Finance	92.2 _{1.0}	94.7 _{1.1}	71.8 _{1.0}	91.9 _{0.4}	76.0 _{0.6}	76.6 _{0.4}	81.9 _{0.6}	78.0 _{0.9}	76.1 _{0.7}	78.3 _{0.6}	81.4 _{0.0}	80.6 _{0.4}	92.3 _{0.0}	72.9 _{8.4}	73.6 _{18.0}	74.3 _{32.5}
SST2 AAVE	85.3 _{0.2}	80.1 _{0.8}	51.9 _{0.1}	81.1 _{0.5}	85.3 _{0.2}	49.2 _{0.3}	49.1 _{0.1}	49.4 _{0.3}	84.3 _{0.2}	86.2 _{0.2}	72.7 _{0.2}	84.9 _{0.2}	64.2 _{0.0}	75.3 _{3.3}	77.6 _{9.1}	72.3 _{11.6}
SST2 SG	83.2 _{0.3}	75.7 _{0.6}	63.7 _{0.5}	81.4 _{0.6}	83.7 _{0.1}	49.3 _{0.3}	49.3 _{0.2}	49.4 _{0.3}	83.6 _{0.1}	85.1 _{0.3}	71.4 _{0.2}	80.6 _{0.4}	64.7 _{0.0}	73.4 _{3.0}	75.3 _{9.7}	70.7 _{10.8}

Table 8. Macro-F1 with GPT2-XL. We report $average_{std}$ over 10 runs for zero-shot. For few-shot, we report $average_{std}$ over 5 seeds.

On Distributional Robustness of In-Context Learning for Text Classification

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	88.7 _{0.3}	89.8 _{0.3}	87.7 _{0.3}	86.8 _{0.4}	75.7 _{0.2}	87.3 _{0.2}	82.1 _{0.3}	84.6 _{0.3}	69.3 _{0.4}	88.4 _{0.2}	80.6 _{0.3}	85.6 _{0.4}	67.6 _{0.0}	83.2 _{2.5}	88.0 _{6.1}	87.8 _{5.6}
SST2 ADV	46.2 _{1.2}	48.9 _{1.1}	60.3 _{1.1}	47.5 _{1.2}	53.2 _{0.8}	64.2 _{1.1}	60.7 _{1.1}	59.5 _{0.9}	51.5 _{0.9}	61.9 _{0.9}	60.8 _{1.2}	61.4 _{0.7}	56.1 _{0.0}	56.4 _{1.7}	54.9 _{1.9}	54.9 _{1.5}
MR	84.8 _{0.2}	85.2 _{0.3}	81.9 _{0.5}	77.7 _{0.4}	76.1 _{0.3}	83.3 _{0.3}	80.6 _{0.3}	80.2 _{0.3}	71.4 _{0.2}	84.1 _{0.3}	80.3 _{0.2}	82.1 _{0.3}	74.7 _{0.0}	83.9 _{1.4}	85.3 _{3.9}	86.7 _{2.5}
CR	87.1 _{0.3}	77.3 _{0.6}	51.0 _{1.6}	80.1 _{0.8}	77.7 _{0.6}	69.5 _{0.5}	67.3 _{0.6}	63.7 _{0.5}	78.5 _{0.8}	70.8 _{0.7}	66.6 _{0.6}	64.9 _{0.6}	78.5 _{0.0}	84.0 _{3.1}	86.9 _{3.0}	86.1 _{3.7}
Amazon	68.6 _{0.3}	91.2 _{0.4}	89.3 _{0.3}	85.1 _{0.3}	75.1 _{0.3}	90.2 _{0.3}	90.1 _{0.2}	89.2 _{0.1}	75.1 _{0.3}	90.5 _{0.2}	90.3 _{0.2}	89.3 _{0.2}	83.9 _{0.0}	91.6 _{0.7}	90.9 _{1.7}	90.6 _{2.3}
Yelp	84.4 _{0.1}	87.7 _{0.3}	73.2 _{0.3}	74.2 _{0.4}	84.4 _{0.1}	87.7 _{0.2}	86.2 _{0.1}	84.4 _{0.2}	84.1 _{0.2}	88.1 _{0.2}	85.9 _{0.1}	84.6 _{0.3}	93.3 _{0.0}	94.0 _{0.5}	92.4 _{1.7}	92.0 _{2.3}
Poem	91.0 _{1.7}	66.0 _{4.8}	69.2 _{1.9}	61.4 _{2.7}	74.6 _{1.9}	70.4 _{1.3}	68.4 _{0.8}	66.2 _{1.1}	74.6 _{1.3}	70.6 _{1.6}	68.4 _{0.8}	65.8 _{1.5}	60.0 _{0.0}	80.6 _{2.8}	65.1 _{11.6}	52.6 _{3.4}
Finance	88.3 _{1.0}	82.0 _{1.4}	90.7 _{1.3}	62.4 _{2.6}	81.0 _{0.6}	89.2 _{0.7}	70.5 _{0.5}	86.8 _{0.6}	81.0 _{0.7}	89.0 _{0.4}	71.4 _{1.0}	86.6 _{0.3}	87.9 _{0.0}	74.1 _{11.7}	85.9 _{6.9}	74.3 _{32.5}
SST2 AAVE	85.0 _{0.3}	69.0 _{0.7}	72.8 _{0.6}	78.4 _{0.5}	85.3 _{0.2}	49.1 _{0.4}	49.0 _{0.4}	49.0 _{0.3}	84.5 _{0.2}	64.7 _{0.5}	83.0 _{0.3}	77.2 _{0.4}	79.0 _{0.0}	83.4 _{3.0}	87.5 _{4.6}	88.2 _{2.6}
SST2 SG	82.4 _{0.3}	62.5 _{0.5}	57.4 _{0.3}	72.9 _{0.6}	83.6 _{0.2}	49.1 _{0.5}	49.0 _{0.3}	49.0 _{0.4}	83.9 _{0.2}	70.4 _{0.4}	83.1 _{0.3}	78.0 _{0.4}	76.4 _{0.0}	81.3 _{3.2}	86.5 _{4.5}	87.3 _{3.3}

Table 9. Macro-F1 with OPT 1.3B. We report $average_{std}$ over 10 runs for zero-shot. For few-shot, we report $average_{std}$ over 5 seeds.

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	89.1 _{0.2}	84.3 _{0.6}	75.3 _{0.7}	82.8 _{0.5}	74.9 _{0.2}	81.3 _{0.3}	77.8 _{0.1}	82.5 _{0.3}	65.3 _{0.4}	79.3 _{0.3}	75.7 _{0.2}	83.6 _{0.3}	57.7 _{0.0}	81.9 _{1.5}	82.3 _{6.0}	86.2 _{5.3}
SST2 ADV	57.0 _{1.0}	46.1 _{1.6}	51.1 _{2.4}	52.4 _{1.1}	51.4 _{0.7}	45.9 _{0.7}	50.5 _{0.7}	52.2 _{0.6}	51.7 _{0.8}	46.1 _{0.9}	48.9 _{0.6}	51.2 _{0.7}	52.0 _{0.0}	43.2 _{3.4}	42.0 _{4.1}	45.5 _{2.2}
MR	81.3 _{0.4}	80.7 _{0.7}	51.9 _{0.5}	81.2 _{0.5}	72.5 _{0.4}	80.3 _{0.2}	76.9 _{0.2}	79.4 _{0.3}	64.6 _{0.4}	78.5 _{0.3}	76.7 _{0.4}	81.4 _{0.2}	78.4 _{0.0}	80.2 _{3.1}	77.9 _{6.2}	80.7 _{9.4}
CR	84.5 _{0.6}	79.2 _{0.6}	37.2 _{0.5}	71.3 _{1.0}	76.2 _{0.5}	77.3 _{0.4}	73.6 _{0.3}	79.6 _{0.3}	76.9 _{0.5}	76.6 _{0.5}	73.5 _{0.2}	79.1 _{0.6}	85.4 _{0.0}	82.6 _{6.8}	83.6 _{3.6}	85.7 _{2.4}
Amazon	88.7 _{0.3}	90.4 _{0.4}	80.5 _{0.4}	89.7 _{0.3}	81.9 _{0.2}	88.0 _{0.2}	84.9 _{0.1}	89.6 _{0.2}	82.1 _{0.2}	87.5 _{0.2}	84.7 _{0.1}	89.4 _{0.2}	86.9 _{0.0}	87.1 _{3.0}	85.7 _{4.9}	88.0 _{2.2}
Yelp	89.1 _{0.2}	91.8 _{0.2}	72.1 _{0.4}	91.9 _{0.2}	87.2 _{0.3}	82.9 _{0.2}	78.1 _{0.1}	84.3 _{0.1}	87.2 _{0.2}	82.8 _{0.2}	77.5 _{0.1}	84.4 _{0.1}	91.6 _{0.0}	88.4 _{4.5}	83.3 _{8.5}	85.8 _{6.1}
Poem	87.0 _{1.7}	79.4 _{2.8}	62.6 _{3.8}	81.2 _{3.9}	79.4 _{1.6}	75.0 _{1.4}	54.0 _{0.0}	75.2 _{1.7}	76.2 _{1.1}	73.0 _{2.2}	54.6 _{1.0}	73.2 _{1.9}	54.3 _{0.0}	81.1 _{6.9}	69.7 _{12.9}	52.6 _{3.4}
Finance	90.4 _{0.9}	86.7 _{1.0}	61.7 _{1.9}	82.9 _{1.0}	86.2 _{0.8}	82.6 _{0.7}	43.4 _{0.6}	86.6 _{0.4}	87.7 _{1.0}	80.9 _{0.5}	46.1 _{0.3}	85.4 _{0.3}	95.3 _{0.0}	92.8 _{1.5}	48.3 _{25.9}	74.3 _{32.5}
SST2 AAVE	85.1 _{0.3}	64.7 _{0.6}	50.0 _{0.1}	82.4 _{0.5}	84.9 _{0.2}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	85.0 _{0.3}	78.5 _{0.5}	53.2 _{0.1}	76.7 _{0.4}	76.3 _{0.0}	81.7 _{1.6}	80.0 _{6.2}	80.5 _{10.0}
SST2 SG	84.5 _{0.4}	74.6 _{0.7}	49.4 _{0.1}	73.5 _{0.6}	84.4 _{0.2}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	85.2 _{0.2}	76.0 _{0.3}	52.8 _{0.2}	77.9 _{0.3}	76.5 _{0.0}	80.9 _{1.1}	78.4 _{6.3}	79.6 _{10.0}

Table 10. Macro-F1 with Pythia 1.4B. We report $average_{std}$ over 10 runs for zero-shot. For few-shot, we report $average_{std}$ over 5 seeds.

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	90.2 _{0.2}	73.2 _{0.5}	90.5 _{0.3}	83.2 _{0.6}	81.6 _{0.4}	76.3 _{0.4}	83.4 _{0.2}	74.5 _{0.4}	78.2 _{0.3}	79.9 _{0.5}	80.0 _{0.4}	83.9 _{0.4}	78.7 _{0.0}	87.6 _{4.5}	92.1 _{2.7}	94.1 _{0.9}
SST2 ADV	53.7 _{1.6}	51.3 _{2.0}	51.1 _{1.6}	50.5 _{1.6}	57.5 _{0.7}	52.4 _{1.1}	66.8 _{0.6}	53.3 _{0.8}	56.4 _{1.1}	55.2 _{1.0}	67.4 _{0.7}	54.9 _{1.7}	52.0 _{0.0}	59.9 _{4.6}	61.9 _{7.3}	63.2 _{4.5}
MR	86.4 _{0.2}	71.4 _{0.6}	76.3 _{0.7}	65.3 _{0.5}	78.2 _{0.3}	71.2 _{0.3}	80.0 _{0.3}	68.8 _{0.3}	74.9 _{0.3}	74.9 _{0.2}	77.2 _{0.4}	78.5 _{0.4}	74.7 _{0.0}	84.5 _{3.7}	84.9 _{6.9}	89.6 _{1.7}
CR	80.3 _{0.4}	63.6 _{1.0}	76.1 _{1.2}	68.6 _{1.1}	78.0 _{0.3}	65.6 _{0.7}	70.5 _{0.4}	67.3 _{0.9}	78.3 _{0.6}	64.0 _{1.0}	66.5 _{0.5}	65.5 _{0.8}	80.6 _{0.0}	85.3 _{2.7}	87.9 _{2.9}	88.1 _{1.3}
Amazon	86.0 _{0.3}	60.7 _{0.4}	77.2 _{0.6}	62.4 _{0.4}	87.4 _{0.2}	82.3 _{0.3}	89.8 _{0.2}	83.5 _{0.2}	87.4 _{0.2}	83.1 _{0.3}	90.8 _{0.2}	84.0 _{0.3}	83.8 _{0.0}	90.8 _{1.2}	88.1 _{4.3}	87.8 _{4.3}
Yelp	89.0 _{0.2}	74.9 _{0.5}	74.6 _{0.4}	65.9 _{0.5}	85.9 _{0.1}	66.6 _{0.3}	86.3 _{0.1}	67.8 _{0.2}	86.1 _{0.2}	68.7 _{0.1}	87.3 _{0.1}	70.0 _{0.2}	90.3 _{0.0}	90.8 _{1.9}	86.8 _{4.5}	86.1 _{5.1}
Poem	88.4 _{2.5}	63.8 _{2.9}	71.4 _{2.3}	71.8 _{2.9}	70.2 _{1.1}	71.4 _{1.0}	91.0 _{1.1}	70.0 _{1.6}	73.0 _{1.1}	72.6 _{1.3}	88.6 _{1.3}	69.6 _{2.1}	82.9 _{0.0}	77.7 _{11.2}	64.6 _{3.9}	52.6 _{3.4}
Finance	91.6 _{0.4}	77.5 _{2.3}	49.3 _{1.2}	64.3 _{2.3}	81.9 _{0.6}	88.4 _{0.7}	75.5 _{1.5}	87.4 _{0.8}	83.2 _{0.7}	88.6 _{0.3}	71.7 _{1.3}	87.3 _{0.3}	90.2 _{0.0}	82.0 _{13.9}	78.7 _{21.0}	74.3 _{32.5}
SST2 AAVE	86.9 _{0.3}	77.4 _{0.6}	60.6 _{0.2}	79.4 _{0.7}	87.0 _{0.3}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	87.3 _{0.4}	71.5 _{0.5}	69.9 _{0.4}	77.5 _{0.4}	79.8 _{0.0}	85.6 _{4.7}	86.9 _{7.2}	92.0 _{1.2}
SST2 SG	85.4 _{0.3}	69.4 _{0.4}	52.4 _{0.3}	74.4 _{0.6}	83.5 _{0.2}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	83.0 _{0.2}	74.6 _{0.3}	73.7 _{0.3}	76.7 _{0.5}	79.7 _{0.0}	84.5 _{5.2}	85.3 _{7.0}	91.4 _{1.6}

Table 11. Macro-F1 with OPT 2.7B. We report $average_{std}$ over 10 runs for zero-shot. For few-shot, we report $average_{std}$ over 5 seeds. Further results are provided in the appendix

On Distributional Robustness of In-Context Learning for Text Classification

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	91.4 _{0.2}	79.9 _{0.6}	88.5 _{0.3}	85.4 _{0.5}	80.1 _{0.4}	81.9 _{0.3}	77.6 _{0.4}	82.8 _{0.5}	71.4 _{0.3}	70.4 _{0.3}	78.1 _{0.5}	69.3 _{0.4}	60.2 _{0.0}	87.3 _{1.4}	90.3 _{1.5}	91.0 _{1.5}
SST2 ADV	47.8 _{0.9}	53.7 _{1.8}	58.9 _{2.0}	51.3 _{1.8}	47.7 _{0.7}	55.7 _{1.0}	55.9 _{1.2}	54.3 _{0.8}	47.7 _{0.8}	57.1 _{1.0}	57.3 _{0.9}	53.4 _{0.9}	45.3 _{0.0}	49.2 _{2.5}	52.7 _{2.1}	52.7 _{1.7}
MR	86.0 _{0.3}	80.3 _{0.5}	72.8 _{0.5}	79.1 _{0.7}	74.9 _{0.3}	76.5 _{0.3}	74.6 _{0.3}	77.6 _{0.3}	67.6 _{0.4}	66.3 _{0.3}	76.5 _{0.3}	64.2 _{0.4}	73.5 _{0.0}	84.6 _{2.1}	87.0 _{0.9}	87.0 _{2.3}
CR	86.1 _{0.4}	78.7 _{0.9}	52.3 _{0.8}	79.0 _{1.1}	85.3 _{0.2}	71.9 _{0.4}	70.7 _{0.4}	60.9 _{0.9}	84.0 _{0.3}	71.8 _{0.4}	71.0 _{0.4}	62.9 _{0.9}	76.9 _{0.0}	85.9 _{3.1}	86.2 _{3.5}	83.3 _{4.2}
Amazon	89.2 _{0.4}	72.7 _{0.6}	80.1 _{0.7}	71.2 _{0.6}	87.4 _{0.3}	86.1 _{0.3}	81.0 _{0.2}	84.0 _{0.2}	87.6 _{0.2}	82.6 _{0.4}	84.1 _{0.2}	81.7 _{0.3}	85.6 _{0.0}	88.3 _{3.0}	91.8 _{2.1}	90.2 _{4.3}
Yelp	87.0 _{0.3}	90.0 _{0.2}	81.2 _{0.2}	85.6 _{0.4}	86.9 _{0.2}	68.1 _{0.2}	66.5 _{0.3}	64.9 _{0.2}	86.8 _{0.2}	67.8 _{0.2}	66.4 _{0.1}	64.4 _{0.2}	87.8 _{0.0}	87.5 _{4.7}	93.4 _{2.3}	85.7 _{9.1}
Poem	84.6 _{1.6}	80.2 _{2.6}	53.2 _{3.0}	73.6 _{2.5}	95.4 _{1.9}	70.0 _{1.6}	71.0 _{1.9}	58.4 _{0.8}	94.4 _{1.8}	69.4 _{1.9}	66.0 _{2.5}	55.0 _{1.4}	60.0 _{0.0}	76.6 _{14.2}	58.3 _{5.0}	52.6 _{3.4}
Finance	93.8 _{0.6}	79.9 _{1.9}	37.8 _{1.3}	73.8 _{1.3}	86.3 _{0.6}	81.2 _{0.6}	68.4 _{0.3}	78.9 _{0.5}	85.5 _{0.6}	78.9 _{0.7}	67.9 _{0.8}	76.9 _{1.0}	58.3 _{0.0}	94.3 _{2.3}	88.2 _{4.4}	74.3 _{32.5}
SST2 AAVE	88.0 _{0.3}	71.9 _{0.9}	59.8 _{0.5}	70.2 _{1.2}	87.6 _{0.1}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	89.1 _{0.3}	71.0 _{0.5}	79.3 _{0.3}	69.8 _{0.4}	73.1 _{0.0}	86.1 _{1.9}	87.3 _{1.5}	87.9 _{2.8}
SST2 SG	87.1 _{0.4}	66.6 _{0.7}	54.5 _{0.2}	61.5 _{0.9}	86.1 _{0.2}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	85.1 _{0.2}	73.0 _{0.5}	73.1 _{0.3}	68.1 _{0.5}	72.3 _{0.0}	84.5 _{2.6}	85.0 _{2.3}	86.3 _{3.6}

Table 12. Macro-F1 with Pythia 2.8B. We report $average_{std}$ over 10 runs for zero-shot. For few-shot, we report $average_{std}$ over 5 seeds.

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	90.7 _{0.2}	86.7 _{0.3}	49.9 _{0.2}	82.0 _{0.5}	80.6 _{0.3}	84.4 _{0.2}	73.2 _{0.2}	83.3 _{0.3}	72.2 _{0.4}	82.8 _{0.5}	68.8 _{0.3}	81.9 _{0.4}	50.7 _{0.0}	88.9 _{3.2}	93.6 _{0.8}	93.6 _{1.3}
SST2 ADV	52.5 _{1.2}	57.1 _{1.7}	55.3 _{1.2}	55.4 _{1.6}	53.4 _{0.6}	53.5 _{1.0}	59.2 _{0.7}	60.3 _{1.2}	53.0 _{0.9}	60.2 _{0.8}	58.3 _{0.5}	56.2 _{1.6}	47.3 _{0.0}	52.4 _{5.6}	58.2 _{3.7}	56.4 _{6.1}
MR	84.7 _{0.3}	75.9 _{0.6}	53.2 _{0.3}	65.9 _{1.0}	77.3 _{0.4}	80.4 _{0.1}	71.8 _{0.3}	80.9 _{0.2}	69.7 _{0.4}	80.5 _{0.2}	67.4 _{0.2}	80.4 _{0.3}	55.3 _{0.0}	85.7 _{2.2}	90.2 _{0.7}	89.5 _{1.4}
CR	86.4 _{0.6}	68.6 _{0.9}	36.9 _{0.2}	60.4 _{1.2}	86.5 _{0.4}	70.6 _{0.5}	68.1 _{0.6}	66.1 _{0.8}	86.0 _{0.2}	67.7 _{0.8}	56.9 _{0.4}	65.9 _{0.5}	68.4 _{0.0}	87.7 _{0.8}	87.3 _{1.8}	86.8 _{3.1}
Amazon	89.3 _{0.3}	82.2 _{0.5}	65.5 _{0.5}	74.2 _{0.5}	87.6 _{0.2}	50.1 _{0.3}	50.1 _{0.3}	50.1 _{0.3}	86.9 _{0.2}	85.5 _{0.2}	89.7 _{0.1}	84.6 _{0.2}	83.8 _{0.0}	93.6 _{0.1}	93.7 _{0.3}	93.5 _{0.2}
Yelp	89.8 _{0.2}	75.3 _{0.3}	66.0 _{0.4}	76.8 _{0.2}	85.9 _{0.1}	50.0 _{0.2}	50.0 _{0.2}	50.0 _{0.2}	85.4 _{0.2}	84.1 _{0.3}	50.2 _{0.2}	82.6 _{0.3}	95.1 _{0.0}	96.2 _{0.5}	95.5 _{0.4}	94.2 _{0.6}
Poem	96.0 _{1.3}	72.0 _{2.5}	62.0 _{3.1}	68.8 _{3.6}	88.8 _{1.0}	78.8 _{1.4}	63.8 _{1.1}	71.0 _{1.1}	88.0 _{1.3}	76.8 _{1.4}	51.4 _{1.6}	75.2 _{1.9}	54.3 _{0.0}	85.7 _{13.3}	80.0 _{9.9}	52.6 _{3.4}
Finance	92.8 _{0.8}	79.3 _{0.7}	71.8 _{1.4}	85.0 _{0.7}	87.7 _{1.0}	83.4 _{1.0}	84.6 _{0.9}	73.3 _{0.5}	87.1 _{0.6}	89.5 _{0.8}	92.9 _{1.2}	78.7 _{0.5}	92.3 _{0.0}	95.1 _{0.6}	90.4 _{5.0}	74.3 _{32.5}
SST2 AAVE	87.0 _{0.4}	70.3 _{0.6}	49.1 _{0.1}	81.2 _{0.2}	87.6 _{0.2}	51.2 _{0.5}	49.5 _{0.2}	55.2 _{0.5}	87.5 _{0.2}	80.9 _{0.4}	49.2 _{0.1}	76.4 _{0.2}	52.9 _{0.0}	88.0 _{2.4}	90.8 _{2.1}	90.4 _{2.8}
SST2 SG	86.5 _{0.2}	90.1 _{0.3}	52.7 _{0.7}	73.2 _{0.5}	86.7 _{0.2}	50.6 _{0.4}	50.1 _{0.2}	49.1 _{0.0}	86.9 _{0.2}	69.2 _{0.3}	53.2 _{0.3}	73.4 _{0.4}	52.5 _{0.0}	86.4 _{2.6}	89.3 _{1.8}	88.9 _{2.9}

Table 13. Macro-F1 with Pythia 6.9B. We report $average_{std}$ over 10 runs for zero-shot. For few-shot, we report $average_{std}$ over 5 seeds.

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	92.1 _{0.3}	82.1 _{0.5}	68.9 _{0.4}	80.3 _{1.1}	90.8 _{0.2}	85.0 _{0.4}	76.2 _{0.4}	85.2 _{0.3}	85.6 _{0.3}	70.2 _{0.4}	57.0 _{0.3}	64.6 _{0.6}	89.3 _{0.0}	61.1 _{10.0}	71.8 _{11.3}	70.6 _{10.5}
SST2 ADV	59.2 _{0.8}	61.6 _{1.3}	56.1 _{1.7}	62.2 _{2.2}	73.9 _{0.6}	58.0 _{0.7}	66.5 _{0.6}	64.4 _{0.8}	74.1 _{0.8}	60.0 _{0.8}	64.2 _{1.0}	65.7 _{0.9}	73.0 _{0.0}	61.0 _{5.2}	62.7 _{5.8}	61.0 _{5.6}
MR	90.2 _{0.2}	74.8 _{0.7}	63.6 _{0.5}	56.1 _{0.7}	86.4 _{0.2}	78.0 _{0.3}	74.9 _{0.3}	80.8 _{0.2}	81.3 _{0.2}	65.1 _{0.3}	57.8 _{0.4}	60.7 _{0.4}	86.9 _{0.0}	56.7 _{6.9}	70.2 _{10.7}	61.6 _{6.9}
CR	86.1 _{0.3}	72.0 _{1.2}	50.8 _{1.5}	73.4 _{1.6}	87.8 _{0.4}	73.9 _{0.7}	76.8 _{1.0}	69.2 _{0.6}	87.5 _{0.4}	76.4 _{0.6}	76.2 _{0.6}	70.7 _{0.8}	51.1 _{0.0}	48.0 _{9.8}	67.8 _{8.1}	71.5 _{8.1}
Amazon	88.4 _{0.4}	92.3 _{0.3}	87.5 _{0.2}	86.2 _{0.3}	91.2 _{0.2}	93.0 _{0.2}	80.1 _{0.3}	89.6 _{0.2}	91.4 _{0.3}	92.3 _{0.1}	79.7 _{0.2}	89.2 _{0.2}	89.7 _{0.0}	59.6 _{10.4}	74.0 _{7.9}	82.9 _{5.1}
Yelp	93.7 _{0.1}	87.2 _{0.2}	90.4 _{0.3}	77.1 _{0.4}	92.4 _{0.1}	92.7 _{0.1}	91.0 _{0.2}	90.0 _{0.2}	92.3 _{0.1}	92.8 _{0.1}	92.2 _{0.1}	89.6 _{0.2}	84.7 _{0.0}	66.1 _{8.6}	72.6 _{9.2}	82.5 _{6.8}
Poem	91.6 _{0.8}	89.2 _{1.0}	84.4 _{1.8}	82.6 _{1.6}	89.6 _{0.8}	82.0 _{0.9}	67.2 _{2.7}	74.0 _{1.6}	90.0 _{0.0}	82.0 _{0.0}	66.0 _{2.7}	73.0 _{1.4}	68.6 _{0.0}	58.3 _{6.4}	59.4 _{6.9}	54.3 _{0.0}
Finance	96.4 _{0.0}	97.1 _{0.6}	46.1 _{2.6}	81.3 _{0.6}	91.3 _{0.9}	81.9 _{0.5}	60.2 _{0.9}	69.8 _{0.5}	90.7 _{0.6}	82.2 _{1.2}	60.7 _{0.6}	69.4 _{1.2}	58.6 _{0.0}	86.5 _{4.4}	73.7 _{32.1}	90.6 _{0.0}
SST2 AAVE	88.8 _{0.3}	65.3 _{0.6}	63.8 _{0.5}	60.2 _{0.4}	89.1 _{0.1}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	89.3 _{0.2}	58.2 _{0.4}	53.1 _{0.3}	57.4 _{0.4}	86.4 _{0.0}	59.7 _{7.8}	71.4 _{10.1}	63.6 _{10.1}
SST2 SG	87.3 _{0.2}	56.4 _{0.4}	50.3 _{0.3}	55.2 _{0.3}	88.6 _{0.2}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	88.6 _{0.3}	62.2 _{0.5}	63.0 _{0.4}	62.2 _{0.4}	84.5 _{0.0}	61.1 _{8.1}	69.3 _{9.1}	65.2 _{11.5}

Table 14. Macro-F1 with Mistral 7B. We report $average_{std}$ over 10 runs for zero-shot. For few-shot, we report $average_{std}$ over 5 seeds. Further results are provided in the appendix

On Distributional Robustness of In-Context Learning for Text Classification

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	90.5 _{0.3}	56.5 _{1.0}	49.7 _{0.8}	53.9 _{0.6}	81.2 _{0.3}	79.8 _{0.4}	52.1 _{0.2}	75.3 _{0.5}	73.1 _{0.3}	73.1 _{0.3}	49.7 _{0.1}	73.1 _{0.4}	68.9 _{0.0}	49.8 _{1.0}	56.4 _{13.4}	60.8 _{13.6}
SST2 ADV	54.7 _{1.7}	57.9 _{1.0}	59.3 _{2.0}	60.4 _{0.9}	52.8 _{0.7}	62.2 _{1.4}	50.1 _{0.3}	64.5 _{1.3}	53.0 _{0.8}	58.8 _{1.4}	49.3 _{0.5}	61.1 _{1.2}	52.7 _{0.0}	54.2 _{3.5}	52.7 _{3.5}	54.3 _{3.3}
MR	88.6 _{0.2}	58.5 _{0.9}	51.6 _{0.3}	60.6 _{1.1}	76.3 _{0.4}	74.2 _{0.4}	52.9 _{0.2}	73.0 _{0.3}	69.5 _{0.3}	68.9 _{0.4}	51.4 _{0.4}	71.5 _{0.2}	67.6 _{0.0}	50.6 _{0.4}	55.5 _{9.2}	56.2 _{8.9}
CR	82.7 _{0.6}	56.9 _{1.6}	60.7 _{1.3}	53.9 _{1.3}	82.6 _{0.3}	63.5 _{0.7}	37.1 _{0.2}	62.4 _{0.5}	82.7 _{0.4}	61.5 _{0.8}	37.2 _{0.2}	58.8 _{0.9}	49.7 _{0.0}	48.5 _{9.3}	51.9 _{8.7}	53.5 _{13.4}
Amazon	79.5 _{0.2}	53.5 _{0.6}	62.5 _{0.6}	55.3 _{0.9}	87.5 _{0.2}	82.0 _{0.4}	57.5 _{0.4}	80.5 _{0.4}	87.6 _{0.2}	82.2 _{0.4}	57.6 _{0.3}	80.1 _{0.3}	69.1 _{0.0}	51.2 _{0.7}	66.6 _{3.9}	62.1 _{7.2}
Yelp	90.4 _{0.3}	59.9 _{0.4}	67.5 _{0.5}	60.9 _{0.3}	90.6 _{0.2}	74.0 _{0.2}	52.9 _{0.2}	69.1 _{0.3}	90.5 _{0.1}	74.0 _{0.3}	52.9 _{0.1}	68.0 _{0.2}	60.1 _{0.0}	50.4 _{0.5}	63.5 _{6.7}	62.2 _{6.0}
Poem	76.4 _{1.6}	53.4 _{6.3}	51.8 _{5.5}	50.8 _{3.4}	85.6 _{1.6}	73.6 _{2.3}	48.0 _{0.0}	64.6 _{1.0}	86.0 _{2.3}	71.0 _{2.4}	48.0 _{0.0}	64.6 _{1.6}	80.0 _{0.0}	51.4 _{3.6}	58.9 _{6.9}	54.3 _{0.0}
Finance	94.1 _{0.4}	88.3 _{3.3}	78.6 _{1.7}	87.7 _{2.2}	93.4 _{0.4}	68.7 _{0.5}	43.8 _{0.3}	69.5 _{0.4}	93.7 _{0.4}	68.0 _{0.7}	44.2 _{0.4}	69.0 _{0.5}	90.9 _{0.0}	73.7 _{32.2}	72.0 _{18.6}	90.6 _{0.0}
SST2 AAVE	82.7 _{0.4}	72.9 _{0.9}	54.9 _{1.1}	69.8 _{0.8}	82.9 _{0.2}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	85.0 _{0.2}	58.3 _{0.4}	49.1 _{0.0}	51.2 _{0.4}	66.6 _{0.0}	50.9 _{1.4}	54.4 _{10.0}	56.3 _{11.7}
SST2 SG	84.1 _{0.3}	56.2 _{1.0}	50.1 _{0.5}	51.2 _{1.2}	85.8 _{0.1}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	86.1 _{0.2}	65.3 _{0.6}	53.7 _{0.3}	70.1 _{0.5}	60.1 _{0.0}	50.9 _{1.4}	54.4 _{9.6}	56.2 _{12.0}

Table 15. Macro-F1 with Llama 7B. We report $average_{std}$ over 10 runs for zero-shot. For few-shot, we report $average_{std}$ over 5 seeds. Further results are provided in the appendix

	Channel simple, ID	Direct context, ID	Direct instruct, ID	Direct simple, ID	Channel simple, ID mix	Direct context, ID mix	Direct instruct, ID mix	Direct simple, ID mix	Channel simple, OOD mix	Direct context, OOD mix	Direct instruct, OOD mix	Direct simple, OOD mix	0-shot, OOD	1-shot, OOD	4-shot, OOD	8-shot, OOD
SST2	92.5 _{0.2}	86.4 _{0.4}	55.8 _{0.4}	85.3 _{0.4}	81.8 _{0.4}	82.6 _{0.6}	59.8 _{0.1}	79.0 _{0.2}	73.5 _{0.3}	81.9 _{0.2}	55.60.1	82.7 _{0.3}	64.7 _{0.0}	56.8 _{15.1}	65.6 _{11.8}	69.0 _{8.9}
SST2 ADV	70.9 _{1.1}	51.2 _{0.8}	50.1 _{1.1}	49.3 _{0.7}	60.8 _{0.5}	58.1 _{0.8}	51.4 _{0.0}	59.7 _{1.3}	60.1 _{0.8}	56.1 _{1.5}	52.0 _{0.0}	59.3 _{1.2}	58.1 _{0.0}	55.5 _{4.8}	56.8 _{4.1}	54.7 _{3.4}
MR	88.8 _{0.1}	80.3 _{0.9}	64.3 _{0.6}	63.9 _{0.6}	77.3 _{0.3}	79.5 _{0.3}	60.6 _{0.2}	75.1 _{0.5}	71.4 _{0.3}	80.3 _{0.2}	55.5 _{0.2}	80.1 _{0.2}	56.9 _{0.0}	56.6 _{12.8}	58.7 _{7.8}	58.7 _{9.3}
CR	90.3 _{0.3}	73.6 _{1.3}	37.0 _{0.4}	69.1 _{0.8}	86.9 _{0.5}	62.7 _{0.8}	39.1 _{0.1}	60.2 _{0.6}	86.6 _{0.4}	59.7 _{0.6}	39.4 _{0.0}	60.1 _{0.5}	61.7 _{0.0}	42.8 _{12.3}	52.8 _{7.5}	62.0 _{10.0}
Amazon	90.7 _{0.1}	89.7 _{0.3}	60.2 _{1.0}	84.5 _{0.4}	88.5 _{0.3}	88.3 _{0.2}	65.8 _{0.2}	83.7 _{0.2}	88.2 _{0.2}	88.5 _{0.2}	65.8 _{0.2}	83.8 _{0.3}	83.0 _{0.0}	56.6 _{1.7}	70.8 _{7.7}	73.6 _{4.7}
Yelp	93.0 _{0.1}	93.5 _{0.3}	72.1 _{0.5}	69.4 _{0.3}	89.3 _{0.1}	64.7 _{0.2}	70.3 _{0.1}	61.0 _{0.2}	89.2 _{0.1}	64.8 _{0.1}	70.1 _{0.1}	60.1 _{0.2}	54.8 _{0.0}	53.0 _{2.9}	64.7 _{15.7}	67.9 _{11.8}
Poem	81.2 _{1.0}	89.0 _{1.7}	79.8 _{2.0}	87.0 _{3.2}	81.2 _{1.0}	77.8 _{0.6}	50.0 _{0.0}	62.2 _{1.1}	81.4 _{1.9}	79.2 _{1.9}	48.4 _{0.8}	63.4 _{1.9}	60.0 _{0.0}	54.9 _{1.1}	53.7 _{6.9}	54.3 _{0.0}
Finance	96.3 _{0.3}	97.2 _{0.6}	80.8 _{0.8}	92.8 _{1.1}	91.6 _{0.6}	73.3 _{0.3}	47.7 _{0.6}	69.7 _{0.6}	91.6 _{0.7}	73.5 _{0.6}	50.0 _{0.6}	69.8 _{0.6}	88.9 _{0.0}	74.8 _{11.5}	63.0 _{30.8}	90.6 _{0.0}
SST2 AAVE	88.5 _{0.4}	79.8 _{0.5}	49.3 _{0.1}	76.3 _{0.6}	89.2 _{0.2}	49.1 _{0.0}	49.1 _{0.0}	49.1 _{0.0}	89.1 _{0.1}	74.9 _{0.1}	71.0 _{0.7}	75.6 _{0.4}	57.2 _{0.0}	56.2 _{13.2}	58.5 _{8.6}	59.4 _{9.7}
SST2 SG	86.5 _{0.2}	90.1 _{0.3}	52.7 _{0.7}	73.2 _{0.5}	86.7 _{0.2}	50.6 _{0.4}	50.1 _{0.2}	49.1 _{0.0}	86.9 _{0.2}	69.2 _{0.3}	53.2 _{0.3}	73.4 _{0.4}	54.8 _{0.0}	50.1 _{3.1}	49.0 _{11.6}	56.1 _{8.5}

Table 16. Macro-F1 with Llama-2-7B. We report $average_{std}$ over 10 runs for zero-shot. For few-shot, we report $average_{std}$ over 5 seeds.

Task (Distribution Shift)	Label Descriptions
Sentiment (Domain)	”This [DOMAIN] leans [POLARITY]:”; DOMAIN $\in \{\text{text, movie review, RottenTomatoes review, customer review, Amazon review, Yelp review, poem, financial news excerpt}\}$, POLARITY $\in \{\text{positive, negative}\}$
Sentiment (Adversarial)	”This misleading text exhibits a [POLARITY] bias:”; POLARITY $\in \{\text{positive, negative}\}$
Sentiment (Dialect)	”This movie review written in [DIALECT] leans [POLARITY]:”; DIALECT $\in \{\text{African American Vernacular English, Colloquial Singapore English (Singlish), Indian English, Irish English, Nigerian English}\}$, POLARITY $\in \{\text{positive, negative}\}$
Emotion (Dialect)	”This tweet written in [DIALECT] emotes [EMOTION]:”; DIALECT $\in \{\text{African American Vernacular English, Colloquial Singapore English (Singlish), Indian English, Irish English, Nigerian English}\}$, EMOTION $\in \{\text{joy, sadness, anger, fear, love, surprise}\}$
News (Dialect)	”This news written in [DIALECT] is about [TOPIC]:”; DIALECT $\in \{\text{African American Vernacular English, Colloquial Singapore English (Singlish), Indian English, Irish English, Nigerian English}\}$, TOPIC $\in \{\text{world, sports, business, technology}\}$

Table 17. Handwritten label description starter templates. DOMAIN=”text” represents missing domain information. We generate label description variations by asking ChatGPT ”Write 9 paraphrases of this sentence as a Python list”

Name	Format
Direct context	"This is a [DISTRIBUTIONAL INFO]. [INPUT] [LABEL DESCRIPTION]"
Direct instruct	"Is this [DISTRIBUTIONAL INFO] [LIST OF LABEL NAMES]? [INPUT] [LABEL DESCRIPTION]"
Direct simple	"[INPUT] [LABEL DESCRIPTION]"
Attribute	Example
DISTRIBUTIONAL INFO	poem
INPUT	with pale blue berries. in these peaceful shades-
LABEL DESCRIPTION	positive
LIST OF LABEL NAMES	negative or positive

Table 18. Different variations of discriminative (direct) models. In all three settings the label descriptions contain contextual information.