

Using Natural Language Explanations to Improve Robustness of In-context Learning

Anonymous ACL submission

Abstract

Recent studies demonstrated that large language models (LLMs) can excel in many tasks via in-context learning (ICL). However, recent works show that ICL-trained models tend to produce inaccurate results when presented with adversarial inputs. In this work, we investigate whether augmenting ICL with natural language explanations (NLEs) improves the robustness of LLMs on adversarial datasets covering natural language inference and paraphrasing identification. We prompt LLMs with a small set of human-generated NLEs to produce further NLEs, yielding more accurate results than both a zero-shot-ICL setting and using only human-generated NLEs. Our results on five popular LLMs (GPT3.5-turbo, LLaMA2, Vicuna, Zephyr, and Mistral) show that our approach yields over 6% improvement over baseline approaches for eight adversarial datasets: HANS, ISCS, NaN, ST, PICD, PISP, ANLI, and PAWS. Furthermore, previous studies have demonstrated that prompt selection strategies significantly enhance ICL on in-distribution test sets. However, our findings reveal that these strategies do not match the efficacy of our approach for robustness evaluations, resulting in an accuracy drop of 8% compared to the proposed approach.

1 Introduction

The landscape of AI has recently undergone a significant transformation with the advent of large language models (LLMs). These models can produce accurate predictions on test data after observing a small number of demonstrations. Remarkably, they can achieve this based on examples provided directly in their inputs, without explicit retraining or fine-tuning – this learning paradigm is referred to as *in-context learning* (ICL, Brown et al., 2020; Rae et al., 2021). However, ICL struggles to execute complex tasks, such as arithmetic, common-sense, and symbolic reasoning (Rae et al., 2021).

To improve the effectiveness of ICL in solving tasks requiring complex reasoning, Wei et al. (2022b) drew inspiration from natural language explanations (NLEs) to introduce a method denoted as the Chain-of-Thought (CoT) prompting. CoT prompting involves prompting a model with a sequence of intermediate steps or reasoning processes to guide it towards generating more accurate answers.¹ In this work, we denote ICL equipped with NLEs as *X-ICL*. Despite its simplicity, X-ICL has advanced the performance of ICL across a broad range of complex reasoning tasks (Wei et al., 2022b; Wang et al., 2023b).

Similarly to supervised learning, ICL tends to be vulnerable to adversarial examples (Wang et al., 2023a). Previous research shows that improving the robustness of fine-tuned models against such adversarial datasets is possible by fine-tuning with task-relevant NLEs (Chen et al., 2022; Ludan et al., 2023). Inspired by this, we hypothesize that incorporating NLEs into ICL could also improve the robustness of LLMs against adversarial examples. To this end, we evaluate the robustness of X-ICL on eight adversarial datasets: HANS, ISCS, NaN, ST, PICD, PISP, ANLI, and PAWS.

Moreover, the effectiveness of X-ICL so far relies on the availability of human-written NLEs (Wei et al., 2022b), which usually require domain-specific knowledge, making them hard to collect. However, the advent of LLMs uncovered a range of possibilities where LLMs can assist human annotators (Bang et al., 2023; Guo et al., 2023). Motivated by this development, we investigate using three LLMs, namely GPT3.5-turbo, LLaMA2, and Vicuna, to generate NLEs for ICL. We then use human annotators to assess the quality of 200 human-written and LLM-generated NLEs. As shown in

¹CoTs and NLEs are similar concepts, as they both describe the reasoning process behind a decision in natural language; as NLEs were introduced before CoTs (Camburu et al., 2018; Hendricks et al., 2018), we use the former term.

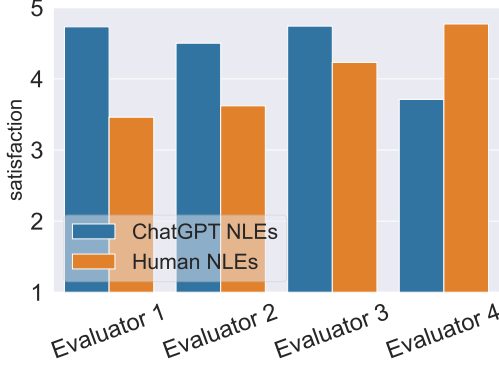


Figure 1: Human evaluation on 100 NLEs generated by GPT3.5-turbo (labeled as *ChatGPT NLEs*) and 100 NLEs generated by human annotators (labeled as *Human NLEs*). The satisfaction scores span from 1 (extremely dissatisfied) to 5 (extremely satisfied).

Figure 1, most annotators (3 out of 4) prefer NLEs produced by ChatGPT (GPT3.5-turbo) over those crafted by humans.² This observation further motivates us to evaluate models trained with LLM-generated NLEs.

We then evaluate the improvement in the robustness of X-ICL in three settings – in two of the settings, an LLM is prompted with LLM-generated NLEs (generated in zero-shot-ICL and few-shots-ICL settings, and in the last setting, the LLM is prompted with human-generated NLEs. In the evaluation, we consider five popular LLMs (*i.e.*, Mistral (Jiang et al., 2023), Zephyr (Tunstall et al., 2023), Vicuna (Chiang et al., 2023), LLaMA2 (Touvron et al., 2023) and GPT3.5-turbo) on eight adversarial datasets. Our experimental results suggest that X-ICL produces more accurate results than ICL and, moreover, that NLEs generated by ChatGPT in a few-shots-ICL setting (by prompting ChatGPT with human-generated NLEs) significantly improve over the ICL baseline (+6%) for the majority of the considered datasets and LLMs. Thus, our findings suggest that an integrated approach, combining human inputs with LLMs, can provide a more effective solution than utilizing either human annotators or LLMs in isolation. Finally, we show that while prompt selection strategies (*i.e.*, retrieving relevant training examples) can significantly improve the accuracy of ICL on in-distribution test sets (Gupta et al., 2023; Levy et al., 2023; Ye et al., 2023), they are less effective on adversarial datasets when compared to X-ICL methods, with our approach (fs-ICL) outperforming them by more than 8% in accuracy.

²More details are available in Appendix D.1.

2 Related Work

Learning with Explanations. There has been a surge of work on explaining predictions of neural NLP systems, from highlighting decision words (Ribeiro et al., 2016; Alvarez-Melis and Jaakkola, 2017; Serrano and Smith, 2019) to generating NLEs (Camburu et al., 2018; Narang et al., 2020; Wiegrefe and Marasovic, 2021). Our work concentrates on the latter category, namely, the self-generation of NLEs for justifying model predictions. Rajani et al. (2019) propose a two-stage training process to improve the prediction performance for commonsense reasoning tasks. In their work, the first stage revolves around generating NLEs, which are then used to inform the label prediction training process in the second stage. Alternatively, one can leverage a multi-task framework to generate NLEs and labels simultaneously (Hase et al., 2020). Li et al. (2022) propose advancing the reasoning abilities of smaller LMs by leveraging NLEs generated by GPT-3 (Brown et al., 2020). NLEs have also vastly been employed beyond NLP, such as in computer vision (Hendricks et al., 2018; Zellers et al., 2019; Majumder et al., 2022), in the medical domain (Kayser et al., 2022), and for self-driving cars (Kim et al., 2018), with some works showing improved task performance when training with NLEs (Kayser et al., 2021). However, these studies primarily concentrate on supervised fine-tuning approaches, which is different from the focus of this work, *i.e.*, ICL.

Prompting with NLEs. Despite its remarkable performance on several downstream tasks (Brown et al., 2020), ICL can still produce inaccurate results in tasks requiring reasoning abilities, such as arithmetic, logical, and commonsense reasoning tasks (Rae et al., 2021; Srivastava et al., 2022). To improve the reasoning abilities of LLMs, Wei et al. (2022b) introduced CoT prompting. This technique prompts a LM to generate a sequence of concise sentences that imitate the reasoning process an individual might undergo to solve a task before providing the ultimate answer, essentially to provide a NLE/CoT before generating the final answer. Furthermore, Wang et al. (2023b) propose to improve CoT prompting by combining multiple diverse reasoning paths generated by LLMs, improving the accuracy of a greedy CoT prompting approach. However, these aforementioned methods need human-written NLEs as CoT in the prompts. Instead, our LLM-based zero-shot-ICL regime har-

nesses the power of an LLM to synthesize NLEs without human-written NLEs.

Learning Robust Models. Several works show that NLP models are prone to performance degradation when presented with adversarial examples, a consequence of inherent artifacts or biases within the annotation of the training dataset (Naik et al., 2018; McCoy et al., 2019; Nie et al., 2020; Liu et al., 2020b). Various strategies have been proposed to mitigate biases within NLP models, *e.g.*, initially training a weak model to recognize superficial features, subsequently enforcing a target model to learn more robust and generalizable characteristics (He et al., 2019; Clark et al., 2019; Karimi Mahabadi et al., 2020; Yaghoobzadeh et al., 2021; Korakakis and Vlachos, 2023). Additionally, data augmentation presents another viable option (Minervini and Riedel, 2018; Wu et al., 2021, 2022). Moreover, studies have shown that supervised fine-tuning of models using rationales or human-written NLEs can significantly enhance the models’ resilience against adversarial datasets (Chen et al., 2022; Stacey et al., 2022; Kavumba et al., 2023; Ludan et al., 2023). Unlike them, our research examines the robustness of X-ICL across eight adversarial datasets, highlighting a novel finding: NLEs generated by LLMs surpass those produced by human annotators in enhancing model robustness. In addition, unlike human-written NLEs, those produced by LLMs exhibit greater scalability and adaptability across diverse tasks.

3 Methodology

This section first outlines the workflow of X-ICL. Then, the focus shifts to detailing how an LLM can generate an NLE for a labeled instance.

3.1 ICL with NLEs (X-ICL)

LLMs can provide significantly more accurate predictions across various reasoning tasks when supplied with human-written NLEs (Wei et al., 2022b,a).

In X-ICL, given an instance, the task is to generate the most likely prediction and NLE for that instance. More formally, in X-ICL, given an unlabeled instance $x' \in \mathcal{X}$ and a set of training examples (x_i, r_i, y_i) , where $x_i \in \mathcal{X}$ is an instance, $y_i \in \mathcal{Y}$ is its label, and $r_i \in \mathcal{E}$ is the corresponding explanation, the task is to identify the most likely

label and explanation for x' :

$$\arg \max_{(r', y') \in \mathcal{E} \times \mathcal{Y}} P_{\theta} \left((r', y') \mid (x_i, r_i, y_i)_{i=1}^k, (x') \right),$$

where θ denotes the model parameters, and \mathcal{X} , \mathcal{Y} , and \mathcal{E} are the sets of all possible instances, labels, and explanations, respectively.

The objective of X-ICL is to maximize the likelihood of generating the optimal NLE, $r' \in \mathcal{E}$, and its corresponding label, $y' \in \mathcal{Y}$, given a demonstration set $(x_i, r_i, y_i)_{i=1}^k$ and an unlabeled instance x' . The most likely combination of label y' and explanation r' is generated by an LLM, after prompting it with the training examples and NLEs $(x_i, r_i, y_i)_{i=1}^k$ and the unlabeled instance x' .

3.2 Generating NLEs with LLMs

In existing X-ICL works, human-written NLEs r were used for the instances within the demonstration set. Instead, in this work, we opt for the NLEs synthesized via LLMs. This preference is driven by noting that NLEs produced by LLMs tend to receive higher approval ratings from human evaluators, as indicated in Figure 1. We argue that this preference will boost the performance of X-ICL. The methods utilized for the generation of NLEs are outlined below.

Few-shot prompting for NLEs Our methodology, also shown in Figure 2, initiates by leveraging a set of labeled instances, each accompanied by a human-crafted NLE, to prompt LLMs. The primary aim is to encourage the LLMs to generate a correct NLE (*i.e.*, the ground-truth arguments) for the correctly predicted answer for a test instance. The most likely NLE is then generated as follows:

$$\arg \max_{r' \in \mathcal{E}} P_{\theta}(r' \mid s, (x_j, y_j, r_j)_{j=1}^m, (x', y')), \quad (1)$$

where s denotes a meta-prompt representing the task. More details on the meta-prompt are available in Appendix B.

Zero-shot prompting for NLEs We further extend our approach to situations where human-written NLEs are absent, which is generally more prevalent across most datasets. In this context, LLMs are prompted to generate an NLE for a labeled instance devoid of any pre-existing examples with NLEs. The objective bears a resemblance to Equation (1), albeit without the inclusion of the demonstration set $(x_j, y_j, r_j)_{j=1}^m$.

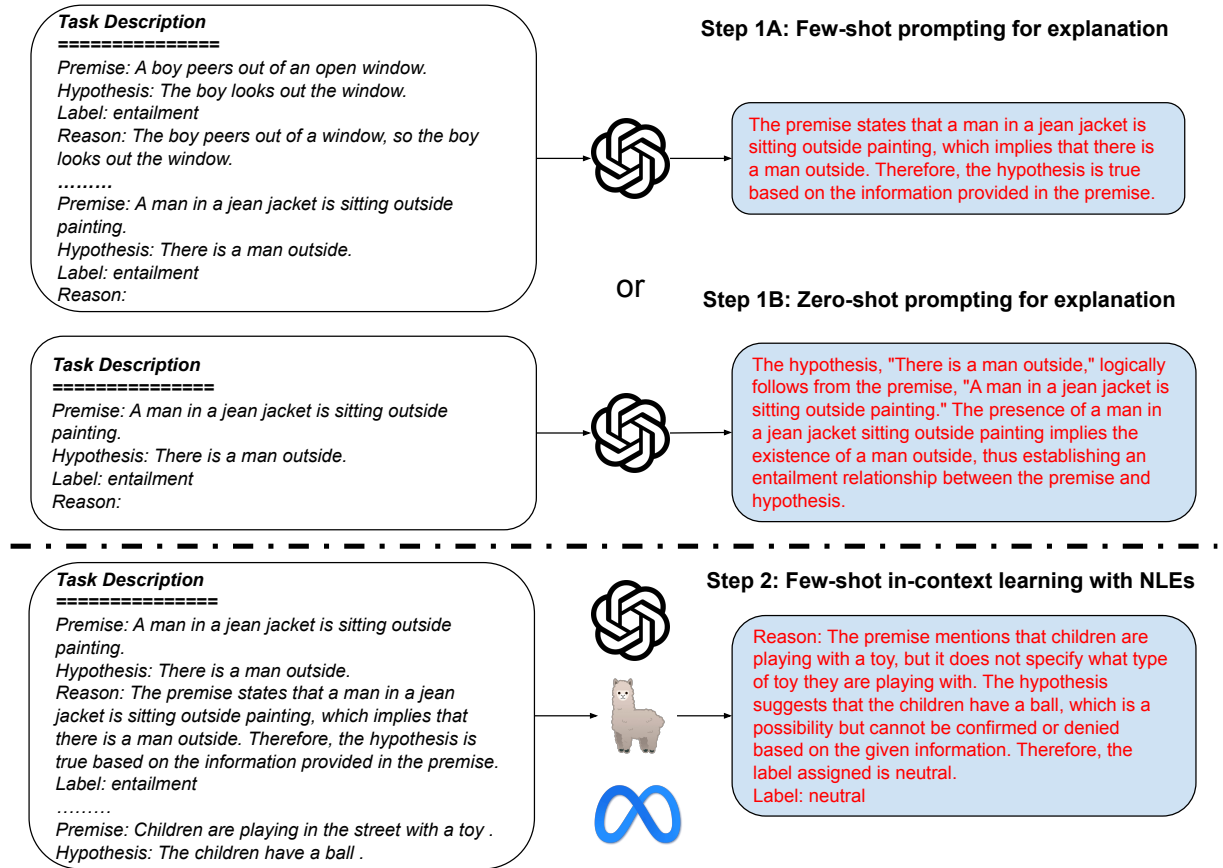


Figure 2: Illustration of using LLM-generated NLEs for ICL: (1) prompt an LLM in a few-shot or zero-shot manner to generate NLEs for new instances; (2) prompt LLMs using ICL with the NLEs generated in step 1.

Notably, the NLEs generated by the aforementioned approaches can be seamlessly integrated into the existing X-ICL framework as delineated in Section 3.1. We primarily focus on using GPT-3.5 (more specifically, GPT3.5-turbo-0613 – we will refer to this model as ChatGPT) to synthesize NLEs. Given that LLMs, such as ChatGPT, may have been trained on datasets incorporating NLEs, it challenges the assumption of genuine zero- or few-shot learning scenarios. To clarify terminology and avoid confusion, we redefine ‘zero-shot learning’ as the absence of demonstration sets, and ‘few-shot ICL’ as learning that utilizes a demonstration set. Thus, we denote the aforementioned two approaches as fs-X-ICL (ChatGPT) and zs-X-ICL (ChatGPT), respectively. In addition, we explore the application of two other widely used open-source LLMs for generating NLEs. Detailed results of these experiments are provided in Appendix C.

4 Experiments

We conduct a series of experiments to assess the performance of our proposed X-ICL framework.

4.1 Experimental Setup

Tasks and datasets We consider the Natural Language Inference (NLI) and paraphrasing identification tasks as our testbed. To ascertain the robustness of LLMs when employing the proposed approach, we evaluate it across eight adversarial datasets. For the NLI task, we include HANS, ISCS, ST, PICD, PISP, NaN, and ANLI. The first five datasets (HANS, ISCS, ST, PICD, PISP) are from Liu et al. (2020b), while NaN and ANLI are sourced from Truong et al. (2022) and Nie et al. (2020), respectively. Regarding the paraphrasing identification task, we use the PAWS-QQP (or PAWS) dataset (Zhang et al., 2019).

Additionally, the SNLI dataset (Bowman et al., 2015) and QQP (Wang et al., 2018), which are non-adversarial, are employed for a comparative purpose. The details of these datasets are provided in Appendix A.

Language models and prompts The evaluation of our approach is undertaken across five prominent LLMs: (1) Mistral, (2) Zephyr, (3) Vicuna, (4) LLaMA2, and (5) GPT3.5-turbo (version 0613). Specifically, the Mistral and Zephyr models have

Models	Methods	Natural Language Inference								Paraphrasing		Avg.
		SNLI	HANS	ISCS	NaN	ST	PICD	PISP	ANLI	QQP	PAWS	
Mistral 7B	ICL	59.8 ±3.4	54.0 ±2.2	51.9 ±1.4	55.0 ±1.3	44.4 ±1.7	58.2 ±2.6	23.0 ±2.6	39.8 ±4.6	69.9 ±1.7	68.3 ±2.7	50.3
	X-ICL (Human)	60.0 ±2.0	56.0 ±2.9	54.7 [▽] ±2.5	58.6 [▽] ±2.9	51.7 [▼] ±4.0	56.9 ±3.3	35.8 [▼] ±6.7	43.9 [▼] ±1.7	69.9 ±0.8	66.4 ±1.5	53.5
	zs-X-ICL (ChatGPT)	56.7 ±6.3	51.8 ±5.1	47.7 ±3.5	55.9 ±5.0	44.9 ±4.8	56.7 ±6.6	25.1 ±8.9	28.8 ±4.4	67.3 ±2.3	64.7 ±3.1	46.4
	fs-X-ICL (ChatGPT)	61.8 ±3.1	58.2[▼] ±2.5	57.2[▼] ±2.2	62.4[▼] ±2.6	55.2[▼] ±1.5	59.2 ±2.7	47.6[▼] ±1.8	46.9[▼] ±2.3	70.3 ±1.1	72.5[▽] ±1.3	57.1
Zephyr 7B	ICL	67.1 ±3.4	71.0 ±1.8	63.4 ±1.2	65.7 ±1.8	60.5 ±1.0	64.8 ±1.5	48.4 ±1.4	47.1 ±1.6	76.9 ±0.4	57.7 ±1.1	59.8
	X-ICL (Human)	72.4 [▼] ±4.3	64.3 ±6.7	58.3 ±5.5	62.0 ±5.3	57.0 ±6.3	60.6 ±9.7	52.0 ±6.7	49.4 ±3.0	75.8 ±1.7	61.4 [▽] ±2.3	59.3
	zs-X-ICL (ChatGPT)	67.2 ±3.9	72.7 ±2.6	60.4 ±5.3	64.0 ±5.2	61.4 ±5.7	64.1 ±5.4	50.8 ±5.2	40.9 ±3.8	74.7 ±1.8	59.1 ±2.4	58.1
	fs-X-ICL (ChatGPT)	74.2[▼] ±3.6	77.4[▼] ±2.2	67.0 ±1.6	67.7 ±2.3	69.3[▼] ±1.5	70.0[▼] ±2.1	65.6[▼] ±2.5	52.1[▽] ±2.8	77.3 ±0.9	61.5[▽] ±1.0	65.5
Vicuna 30B	ICL	65.2 ±2.7	69.4 ±1.2	62.7 ±0.9	61.4 ±3.5	58.7 ±0.8	67.1 ±1.6	50.9 ±1.3	50.0 ±2.6	81.8 ±0.5	69.7 ±2.6	61.4
	X-ICL (Human)	67.8 ±3.2	62.9 ±3.7	60.9 ±2.2	64.2 ±1.2	57.3 ±2.0	63.7 ±7.2	55.0 ±5.8	48.2 ±4.7	77.4 ±2.8	63.4 ±3.5	59.8
	zs-X-ICL (ChatGPT)	64.2 ±5.9	61.4 ±7.7	64.9 ±2.3	60.2 ±4.0	61.7 ±3.1	57.9 ±8.7	51.8 ±8.7	49.7 ±3.6	72.1 ±3.2	61.8 ±4.9	58.8
	fs-X-ICL (ChatGPT)	65.0 ±3.1	74.5[▽] ±4.4	65.5[▽] ±1.6	66.3[▽] ±1.1	64.8[▼] ±1.8	61.6 ±8.9	65.9[▼] ±4.7	57.5[▼] ±1.3	78.6 ±1.7	70.0 ±3.3	65.4
LLaMA2 70B	ICL	69.3 ±1.2	65.7 ±3.4	63.1 ±1.6	61.5 ±2.3	58.8 ±4.4	67.6 ±3.0	48.5 ±7.3	54.2 ±2.9	80.8 ±0.6	44.5 ±2.9	60.3
	X-ICL (Human)	73.0 [▼] ±3.1	65.2 ±4.6	59.6 ±4.4	62.4 ±3.3	55.7 ±3.9	64.3 ±2.3	50.4 ±5.1	49.0 ±2.6	74.5 ±3.0	42.6 ±3.3	57.7
	zs-X-ICL (ChatGPT)	55.4 ±5.5	64.0 ±6.3	37.4 ±6.0	58.1 ±5.4	47.7 ±5.4	53.5 ±8.5	44.2 ±8.7	35.8 ±0.8	69.1 ±4.1	37.8 ±4.8	48.1
	fs-X-ICL (ChatGPT)	74.2[▼] ±2.5	73.3[▼] ±8.5	57.7 ±1.2	65.9[▽] ±3.2	63.1[▽] ±3.7	70.6[▽] ±6.5	55.8[▼] ±5.9	59.2[▼] ±1.6	77.6 ±0.6	46.5[▽] ±1.9	63.6
GPT3.5-turbo	ICL	71.9 ±1.4	72.4 ±0.6	64.4 ±0.9	70.0 ±0.8	62.1 ±1.6	64.0 ±3.1	51.2 ±0.4	56.1 ±2.0	81.5 ±0.3	42.9 ±2.8	62.4
	X-ICL (Human)	78.0[▼] ±1.7	71.0 ±1.7	69.0 [▽] ±1.2	70.5 ±2.2	65.7 [▽] ±1.0	72.7 [▼] ±1.3	59.3 [▽] ±1.9	59.8 [▽] ±2.3	76.0 ±3.9	53.4 [▼] ±5.3	66.2
	zs-X-ICL (ChatGPT)	71.9 ±2.7	71.6 ±0.8	68.4 [▽] ±0.3	70.2 ±0.0	67.6 [▽] ±1.3	67.7 [▽] ±4.1	61.7 [▼] ±1.9	60.4[▼] ±2.0	80.4 ±0.8	51.2 [▼] ±3.1	66.0
	fs-X-ICL (ChatGPT)	75.5 [▽] ±2.8	76.0[▼] ±2.0	74.9[▼] ±0.1	73.1[▼] ±1.4	73.3[▼] ±0.4	76.9[▼] ±0.4	75.5[▼] ±3.0	59.6 [▽] ±1.8	79.0 ±1.7	54.0[▼] ±2.6	69.7

Table 1: Accuracy of multiple LLMs using (1) standard ICL without NLEs, (2) X-ICL with human-written NLEs: X-ICL (Human), (3) X-ICL with ChatGPT-generated NLEs in a zero-shot scenario: zs-X-ICL (ChatGPT), (4) X-ICL with ChatGPT-generated NLEs in a few-shot scenario: fs-X-ICL (ChatGPT). The best performance for each task within a model is shown in **bold**. Significance testing was assessed via an unequal variances t -test in comparison with ICL: \blacktriangledown (resp. ∇) represents a p -value lower than 10^{-3} (resp. 10^{-1}). The results of ANLI are the average of ANLI R1, R2, and R3.

7B parameters each. For Vicuna and LLaMA2, we use the 30B and 70B-chat versions, respectively.

We perform all experiments in an 8-shot setting, wherein each experiment is conducted four times independently, thereby drawing 32 unique instances from the training-associated datasets as follows. Specifically, for NLI datasets (except ANLI, which includes its own training set and NLEs) we adhere to the established methodology of using the e-SNLI dataset as the demonstration set, as suggested by Liu et al. (2020b). The e-SNLI dataset is a modified version of SNLI, where each instance is annotated with NLEs written by humans. In the case of the QQP and PAWS datasets, the QQP dataset is utilized as the demonstration set. As no

NLEs are available, we contribute the corresponding NLEs (refer to Appendix F).

Regarding the generation of NLEs via few-shot learning described in section 3.2, the methodology involves selecting a random instance from each label category within the training dataset to form the demonstration set. Consequently, the demonstration set comprises three instances for the e-SNLI dataset and two for the QQP dataset.

Baselines In addition to the proposed method, our study investigates two baselines for comparative analysis. The first baseline uses standard ICL without NLEs. The second employs human-written NLEs within the X-ICL process, referred to as X-ICL (Human).

4.2 Main Results

This section examines ICL and X-ICL across the studied datasets using Mistral, Zephyr, Vicuna, LLaMA2, and GPT3.5-turbo. The results are summarized in Table 1.

The results demonstrate a consistent outcome across both scenarios: with and without the application of X-ICL. As the capabilities of the models increase, there is a noticeable improvement in average accuracy. This progression is evident when comparing the least potent model, exemplified by Mistral, to the most advanced one, represented by GPT3.5-turbo.

Table 1 demonstrates that X-ICL (Human) yields a better predictive accuracy than ICL across all five LLMs assessed using the SNLI dataset, with enhancements of up to 6.1%. This performance elevation is, however, limited to the Mistral and GPT-3.5-turbo models when subjected to all adversarial NLI test sets. The advantage of X-ICL (Human) relative to ICL diminishes when applied to the QQP and PAWS datasets.

For fs-X-ICL (ChatGPT), both Mistral and Zephyr platforms demonstrate a significant performance advantage in all evaluated tasks, outperforming ICL and X-ICL (Human) by at least 5.7% and 3.6%, respectively. Despite the notable improvement on ICL when employing GPT3.5-turbo in comparison to other LLMs, fs-X-ICL (ChatGPT) offers substantially additional gains, with an increase in absolute accuracy between 11%-24% on tasks such as ISCS, ST, PICD, PISP, and PAWS. This suggests that X-ICL enhances LLM effectiveness on in-distribution test sets and increases their robustness against adversarial test sets.

Remarkably, despite the predominant preference of human evaluators for NLEs generated by GPT3.5 over those written by humans, zs-X-ICL (ChatGPT) consistently produces less accurate results than X-ICL (Human) across all models under study. The exception to this trend is GPT3.5-turbo, where a tie is observed. Furthermore, it appears counter-intuitive that zs-X-ICL (ChatGPT) is outperformed by ICL for 4 out of the 5 LLMs analyzed, especially on LLaMA2. We conduct a systematic analysis in section 4.4 to understand this apparent discrepancy between human preferences and LLM performance.

In light of the encompassment of diverse robustness scenarios by the seven adversarial NLI datasets, our primary focus henceforth will be the

Models	Methods	SNLI	AdvNLI	Δ
Zephyr	ICL	67.1	57.2	9.9
	fs-X-ICL (ChatGPT)	74.2	63.7	10.5
	COSINE	77.0	55.6	21.4
	BM25	70.1	53.7	16.4
	SET-BSR	79.9	59.7	20.2
GPT3.5-turbo	ICL	71.9	61.4	10.5
	fs-X-ICL (ChatGPT)	75.5	69.8	5.6
	COSINE	75.0	58.1	16.9
	BM25	71.4	56.0	15.4
	SET-BSR	77.4	59.5	17.9

Table 2: Performance of ICL, fs-X-ICL (ChatGPT) and three data selection approaches on SNLI and AdvNLI (*i.e.*, seven adversarial test sets). Δ indicates the difference between SNLI and adversarial NLI test sets. We report the average performance over all adversarial test sets.

examination of these NLI datasets.

4.3 Impacts of NLEs

Our research has demonstrated that using NLEs generated by GPT3.5 can substantially enhance the performance of X-ICL. To provide a more comprehensive understanding of the NLEs’ influence, we conducted two investigations, presented below.

Data selection vs. X-ICL. The effectiveness of ICL in LLMs is closely linked to the quality of demonstrations provided, as these demonstrations are critical for the model’s ability to understand and address the test instances (Zhao et al., 2021; Liu et al., 2022; Lu et al., 2022). Consequently, considerable research has focused on developing data selection techniques to optimize the curation of ICL demonstrations from relevant candidate data pools, aiming to enhance their alignment with the test instances (Gupta et al., 2023; Levy et al., 2023; Ye et al., 2023). While these approaches have proven to be highly effective on in-distribution test sets, their performance on adversarial test sets remains uncertain, as these sets have the potential to misguide the selection algorithms.

In this context, we compare the performance of fs-X-ICL (ChatGPT) to three prevalent data selection techniques: COSINE, BM25, and SET-BSR. COSINE incorporates sentence embeddings (Reimers and Gurevych, 2019) to identify the most relevant demonstrations for each test instance, while BM25 employs the BM25 algorithm (Sparck Jones et al., 2000) for retrieving candidate demonstrations. SET-BSR utilizes BERTScore (Zhang et al., 2020), integrated with set theory, to ensure comprehensive information coverage and diversity

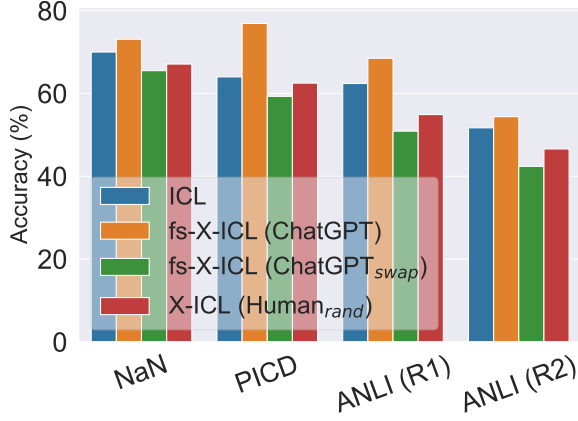


Figure 3: ICL performance of GPT3.5-turbo using (1) standard ICL without NLEs, (2) X-ICL with GPT3.5-generated NLEs in a few-shot scenario: fs-X-ICL (ChatGPT), (3) X-ICL with GPT3.5-generated NLEs, where the NLEs of the prompt are swapped and do not match the instances: fs-X-ICL (ChatGPT_{swap}), and (4) X-ICL with random human NLEs: X-ICL (Human_{rand}).

within the selected instances (Gupta et al., 2023). Note that these data selection techniques are designed to sift through the entirety of the training data to choose demonstrations, a computationally demanding and computationally expensive process for generating NLEs for the full dataset. Therefore, our analysis is confined to applying ICL to these methods. To facilitate a generic comparison with the in-distribution set, we consider the average performance across all adversarial NLI test sets.

According to Table 2, as expected, the data selection approaches markedly enhance ICL performance on the SNLI dataset for all studied LLMs, with notable improvements observed in SET-BSR, achieving gains of up to 17.8% over standard ICL. However, this pronounced advantage diminishes considerably on adversarial test sets, particularly for COSINE and BM25 models, which are outperformed by ICL across all tested LLMs. This discrepancy results in a marked disparity between the in-distribution test set and adversarial test sets, contrary to what is observed in fs-X-ICL (ChatGPT). These results imply that current data selection approaches may be prone to overfitting on in-distribution tests, potentially leading to significant challenges in processing OOD and adversarial datasets due to their limited generalizability.

Do proper NLEs really help? The prevailing assumption argues that the benefits of the X-ICL primarily originate from the NLEs provided. To conclusively attribute these gains to the NLEs rather than any potential influence of additional sentences,

Premise: None of them supported her.

Hypothesis: One of them supported her.

NLE [X-ICL (Human)]: If none of them supported her, then one of them did not support her.

NLE [fs-X-ICL (ChatGPT)]: The hypothesis contradicts the given premise, which states that none of them supported her.

Premise: Not all people have had the opportunities you have had.

Hypothesis: Some people have not had the opportunities you have had.

NLE [X-ICL (Human)]: If not all people have had the opportunities you have had, then some people have not had the opportunities you have had.

NLE [fs-X-ICL (ChatGPT)]: The hypothesis is a direct result of the premise, and the label assigned is entailment.

Table 3: Two test examples from the NaN dataset and the corresponding NLEs generated by X-ICL (Human) and fs-X-ICL (ChatGPT) using Zephyr.

we investigate two experimental setups. In the first setup, we randomly swap the NLEs within the prompt, leading to a mismatched NLE for each instance. This variant is henceforth referred to as fs-X-ICL (ChatGPT_{swap}). Regarding the second variant, for each instance in the demonstration set, we randomly select an unrelated human NLE from the corresponding training set, referred to as X-ICL (Human_{rand}).

As depicted in Figure 3, despite identical content being provided to GPT3.5-turbo, a misalignment between the NLE and the instance results in a marked reduction in the performance of fs-X-ICL (ChatGPT_{swap}) when compared to fs-X-ICL (ChatGPT). This decline is discernible across various datasets, including NaN, PICD, and ANLI (R1/R2).³ It is also shown that an irrelevant and arbitrary NLE triggers a performance reduction within the X-ICL framework. Furthermore, the efficiency of both fs-X-ICL (ChatGPT_{swap}) and X-ICL (Human_{rand}) substantially lags behind that of ICL. Therefore, it can be inferred that the efficacy of the fs-X-ICL (ChatGPT) hinges on providing an accurate and relevant NLE.

4.4 Supplementary Studies

Why is fs-X-ICL (ChatGPT) producing the most accurate results? Our study demonstrates that fs-X-ICL (ChatGPT) surpasses both X-ICL (Human)

³Similar patterns have been detected in other datasets

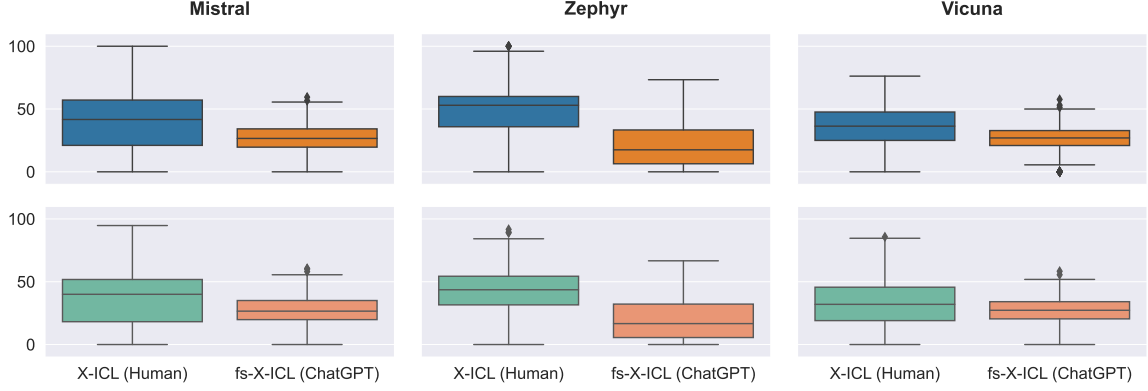


Figure 4: ROUGE-L between the NAN test set and the corresponding generated NLEs. **Top:** ROUGE-L between test premise and NLE. **Bottom:** ROUGE-L between test hypothesis and NLE.

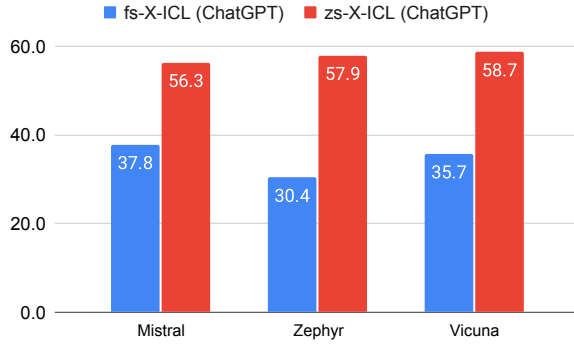


Figure 5: Average length (#words) of NLEs generated by fs-X-ICL (ChatGPT) and zs-X-ICL (ChatGPT).

and zs-X-ICL (ChatGPT) in accuracy. However, the reasons behind this superior performance are not yet understood. Therefore, this section focuses on systematically analyzing the efficacy of fs-X-ICL (ChatGPT).

We first dissect the effectiveness of fs-X-ICL (ChatGPT) over X-ICL (Human). As shown in Table 3, NLEs from X-ICL (Human) are mere verbatim copies of inputs rather than insightful explanations. To substantiate this, we calculate the ROUGE-L scores between the NAN test set and the corresponding NLEs from X-ICL (Human) and fs-X-ICL (ChatGPT) as a means of similarity measurement. As depicted in Figure 4, NLEs from X-ICL (Human) often replicate the given premise and hypothesis, resulting in high ROUGE-L scores. Instead, fs-X-ICL (ChatGPT) can produce meaningful NLEs, demonstrating lower similarity to the test instances.

After analyzing the NLEs from zs-X-ICL (ChatGPT), we attribute the inefficiency to verbose NLEs. Specifically, Figure 5 shows that zs-X-ICL (ChatGPT) produces longer NLEs than fs-X-ICL (ChatGPT). As a result, we observe inconsistency within the NLEs, leading to incorrect pre-

Methods	Mistral	Zephyr	Vicuna
X-ICL (Human)	53.5	59.3	59.8
zs-X-ICL (ChatGPT)	46.4	58.1	58.8
zs-X-ICL (ChatGPT _s)	56.2	62.3	63.4
fs-X-ICL (ChatGPT)	57.1	65.5	62.1

Table 4: Average accuracy of X-ICL (Human), zs-X-ICL (ChatGPT), zs-X-ICL (ChatGPT_s) and fs-X-ICL (ChatGPT) among all test sets.

dictions. As a remedy, we prompt ChatGPT to generate shorter NLEs in the zero-shot setting, denoted as zs-X-ICL (ChatGPT_s). Compared to zs-X-ICL (ChatGPT), the NLEs generated by zs-X-ICL (ChatGPT_s) are reduced to an average of 27 tokens. Consequently, with the help of the concise NLEs, we can improve the accuracy significantly and even surpass the X-ICL (Human) as shown in Table 4.

5 Summary and Outlook

We introduced a simple yet effective method called fs-X-ICL (ChatGPT), leveraging human-written NLEs to generate synthetic NLEs by prompting ChatGPT. fs-X-ICL (ChatGPT) significantly boosts accuracy across various adversarial datasets and five LLMs, compared to standard in-context learning and X-ICL using human-written NLEs. Additionally, our analysis revealed that data selection methodologies may exhibit overfitting within the in-distribution dataset, thus potentially failing to extend to unseen or adversarial datasets. In contrast, our approach employing NLEs has shown consistent performance in both in-distribution and adversarial contexts. Our work paves the way for more robust performance and enhanced explainability capabilities of LLMs.

Limitations

One limitation of X-ICL might be the observed lack of fidelity in the NLEs generated by LLMs, despite their capability to provide accurate answers. These NLEs may sometimes include unfaithful or hallucinated information, which if relied upon by users for model trust, can lead to severe implications. Testing and enhancing the faithfulness of NLEs is a challenging open question (Atanasova et al., 2023). In this work, we show that X-ICL improves robustness, but we do not advocate for the usage of the generated NLEs as faithful explanations without further testing. Second, our approach exhibited promising results when tested against adversarial datasets in two notable NLP tasks: natural language inference and paraphrasing identification. However, further research is required to examine the performance of LLMs and their generalizability across diverse NLP tasks in the context of adversarial examples.

References

- David Alvarez-Melis and Tommi Jaakkola. 2017. [A causal framework for explaining the predictions of black-box sequence-to-sequence models](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 412–421, Copenhagen, Denmark. Association for Computational Linguistics.
- Pepa Atanasova, Oana-Maria Camburu, Christina Lioma, Thomas Lukasiewicz, Jakob Grue Simonsen, and Isabelle Augenstein. 2023. Faithfulness Tests for Natural Language Explanations. In *ACL*.
- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *CoRR*, abs/2302.04023.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *EMNLP*, pages 632–642. The Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Oana-Maria Camburu, Tim Rocktäschel, Thomas Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natural language inference with natural language explanations. *Advances in Neural Information Processing Systems*, 31.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. 2023. [Quantifying memorization across neural language models](#). In *The Eleventh International Conference on Learning Representations*.
- Howard Chen, Jacqueline He, Karthik Narasimhan, and Danqi Chen. 2022. Can rationalization improve robustness? In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3792–3805.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. [Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality](#).
- Christopher Clark, Mark Yatskar, and Luke Zettlemoyer. 2019. [Don’t take the easy way out: Ensemble based methods for avoiding known dataset biases](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4069–4082, Hong Kong, China. Association for Computational Linguistics.
- Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How close is chatgpt to human experts? comparison corpus, evaluation, and detection. *CoRR*, abs/2301.07597.
- Shivanshu Gupta, Sameer Singh, and Matt Gardner. 2023. Coverage-based example selection for in-context learning. *arXiv preprint arXiv:2305.14907*.
- Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. 2018. [Annotation artifacts in natural language inference data](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112, New Orleans, Louisiana. Association for Computational Linguistics.
- Peter Hase, Shiyue Zhang, Harry Xie, and Mohit Bansal. 2020. [Leakage-adjusted simulatability: Can models generate non-trivial explanations of their behavior in natural language?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages

749	Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3428–3448, Florence, Italy. Association for Computational Linguistics.	806
750		807
751		808
752		
753		809
754		810
		811
755	Pasquale Minervini and Sebastian Riedel. 2018. Adversarially regularising neural NLI models to integrate logical background knowledge . In <i>Proceedings of the 22nd Conference on Computational Natural Language Learning</i> , pages 65–74, Brussels, Belgium. Association for Computational Linguistics.	812
756		813
757		
758		814
759		815
760		816
		817
761	Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. Stress test evaluation for natural language inference . In <i>Proceedings of the 27th International Conference on Computational Linguistics</i> , pages 2340–2353, Santa Fe, New Mexico, USA. Association for Computational Linguistics.	818
762		
763		819
764		820
765		821
766		822
767		
768	Sharan Narang, Colin Raffel, Katherine Lee, Adam Roberts, Noah Fiedel, and Karishma Malkan. 2020. Wt5?! training text-to-text models to explain their predictions. <i>arXiv preprint arXiv:2004.14546</i> .	823
769		824
770		825
771		826
		827
772	Yixin Nie, Yicheng Wang, and Mohit Bansal. 2019. Analyzing compositionality-sensitivity of nli models. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 33, pages 6867–6874.	828
773		829
774		
775		830
776	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. Adversarial NLI: A new benchmark for natural language understanding. In <i>ACL</i> , pages 4885–4901. Association for Computational Linguistics.	831
777		832
778		833
779		834
780		
781	Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. 2021. Scaling language models: Methods, analysis & insights from training gopher. <i>arXiv preprint arXiv:2112.11446</i> .	835
782		836
783		837
784		838
785		839
786		840
787	Nazneen Fatema Rajani, Bryan McCann, Caiming Xiong, and Richard Socher. 2019. Explain yourself! leveraging language models for commonsense reasoning . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4932–4942, Florence, Italy. Association for Computational Linguistics.	841
788		842
789		843
790		844
791		845
792		846
793		847
794	Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	848
795		849
796		850
797		851
798		852
799		853
800		854
801		855
802		856
803	Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?": Explaining the predictions of any classifier. In <i>Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16</i> , page 1135–1144, New York, NY, USA. Association for Computing Machinery.	857
804		858
805		859
		860
		861
	Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, and Eneko Agirre. 2023. Did chatgpt cheat on your test?	
	William A. Scott. 1962. Cognitive complexity and cognitive flexibility . <i>Sociometry</i> , 25(4):405–414.	
	Sofia Serrano and Noah A. Smith. 2019. Is attention interpretable? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2931–2951, Florence, Italy. Association for Computational Linguistics.	
	K. Sparck Jones, S. Walker, and S.E. Robertson. 2000. A probabilistic model of information retrieval: development and comparative experiments: Part 1 . <i>Information Processing and Management</i> , 36(6):779–808.	
	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <i>arXiv preprint arXiv:2206.04615</i> .	
	Joe Stacey, Yonatan Belinkov, and Marek Rei. 2022. Supervising model attention with human explanations for robust natural language inference. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pages 11349–11357.	
	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	
	Thinh Hung Truong, Yulia Otmakhova, Timothy Baldwin, Trevor Cohn, Jey Han Lau, and Karin Verspoor. 2022. Not another negation benchmark: The NaN-NLI test suite for sub-clausal negation . In <i>Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 883–894, Online only. Association for Computational Linguistics.	
	Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. 2023. Zephyr: Direct distillation of lm alignment. <i>arXiv preprint arXiv:2310.16944</i> .	
	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding . In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing</i>	

862	<i>and Interpreting Neural Networks for NLP</i> , pages	Yadollah Yaghoobzadeh, Soroush Mehri, Remi Ta-	918
863	353–355, Brussels, Belgium. Association for Com-	chet des Combes, T. J. Hazen, and Alessandro Sor-	919
864	putational Linguistics.	doni. 2021. Increasing robustness to spurious corre-	920
865	Jiong Xiao Wang, Zichen Liu, Keun Hee Park, Muhao	lations using forgettable examples . In <i>Proceedings</i>	921
866	Chen, and Chaowei Xiao. 2023a. Adversarial demon-	<i>of the 16th Conference of the European Chapter of</i>	922
867	stration attacks on large language models. <i>arXiv</i>	<i>the Association for Computational Linguistics: Main</i>	923
868	<i>preprint arXiv:2305.14950</i> .	<i>Volume</i> , pages 3319–3332, Online. Association for	924
869	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le,	Computational Linguistics.	925
870	Ed H. Chi, Sharan Narang, Aakanksha Chowdhery,	Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and	926
871	and Denny Zhou. 2023b. Self-consistency improves	Lingpeng Kong. 2023. Compositional exemplars for	927
872	chain of thought reasoning in language models . In	in-context learning. In <i>Proceedings of the 40th Inter-</i>	928
873	<i>The Eleventh International Conference on Learning</i>	<i>national Conference on Machine Learning, ICML’23</i> .	929
874	<i>Representations</i> .	JMLR.org.	930
875	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,	Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin	931
876	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,	Choi. 2019. From recognition to cognition: Vi-	932
877	Maarten Bosma, Denny Zhou, Donald Metzler, Ed H.	sual commonsense reasoning. In <i>Proceedings of the</i>	933
878	Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy	<i>IEEE/CVF Conference on Computer Vision and Pat-</i>	934
879	Liang, Jeff Dean, and William Fedus. 2022a. Emer-	<i>tern Recognition</i> .	935
880	gent abilities of large language models . <i>Transactions</i>	Tianyi Zhang, Varsha Kishore*, Felix Wu*, Kilian Q.	936
881	<i>on Machine Learning Research</i> . Survey Certifica-	Weinberger, and Yoav Artzi. 2020. Bertscore: Eval-	937
882	tion.	uating text generation with bert . In <i>International</i>	938
883	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	<i>Conference on Learning Representations</i> .	939
884	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,	Yuan Zhang, Jason Baldridge, and Luheng He. 2019.	940
885	and Denny Zhou. 2022b. Chain-of-thought prompt-	PAWS: Paraphrase adversaries from word scrambling .	941
886	ing elicits reasoning in large language models. In	In <i>Proceedings of the 2019 Conference of the North</i>	942
887	<i>NeurIPS</i> .	<i>American Chapter of the Association for Computa-</i>	943
888	Sarah Wiegrefe and Ana Marasovic. 2021. Teach me to	<i>tional Linguistics: Human Language Technologies,</i>	944
889	explain: A review of datasets for explainable natural	<i>Volume 1 (Long and Short Papers)</i> , pages 1298–1308,	945
890	language processing. <i>35th Conference on Neural</i>	Minneapolis, Minnesota. Association for Computa-	946
891	<i>Information Processing Systems (NeurIPS) Track on</i>	tional Linguistics.	947
892	<i>Datasets and Benchmarks</i> .	Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and	948
893	Adina Williams, Nikita Nangia, and Samuel Bowman.	Sameer Singh. 2021. Calibrate before use: Improv-	949
894	2018. A broad-coverage challenge corpus for sen-	ing few-shot performance of language models . In	950
895	tence understanding through inference . In <i>Proceed-</i>	<i>Proceedings of the 38th International Conference</i>	951
896	<i>ings of the 2018 Conference of the North American</i>	<i>on Machine Learning</i> , volume 139 of <i>Proceedings</i>	952
897	<i>Chapter of the Association for Computational Lin-</i>	<i>of Machine Learning Research</i> , pages 12697–12706.	953
898	<i>guistics: Human Language Technologies, Volume</i>	PMLR.	954
899	<i>1 (Long Papers)</i> , pages 1112–1122, New Orleans,		
900	Louisiana. Association for Computational Linguis-		
901	tics.		
902	Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and		
903	Daniel Weld. 2021. Polyjuice: Generating counter-		
904	factuals for explaining, evaluating, and improving		
905	models . In <i>Proceedings of the 59th Annual Meet-</i>		
906	<i>ing of the Association for Computational Linguistics</i>		
907	<i>and the 11th International Joint Conference on Natu-</i>		
908	<i>ral Language Processing (Volume 1: Long Papers)</i> ,		
909	pages 6707–6723, Online. Association for Computa-		
910	tional Linguistics.		
911	Yuxiang Wu, Matt Gardner, Pontus Stenetorp, and		
912	Pradeep Dasigi. 2022. Generating data to mitigate		
913	spurious correlations in natural language inference		
914	datasets . In <i>Proceedings of the 60th Annual Meet-</i>		
915	<i>ing of the Association for Computational Linguistics</i>		
916	<i>(Volume 1: Long Papers)</i> , pages 2660–2676, Dublin,		
917	Ireland. Association for Computational Linguistics.		

A Details of Datasets

The details of all studied datasets are delineated as follows

- **SNLI Dataset:** The SNLI dataset, a benchmark in natural language inference, encompasses approximately 570,000 human-annotated sentence pairs, each pair formed by a premise and a hypothesis. These sentences originate from an existing corpus of image captions, thus offering a broad spectrum of common subjects and linguistic structures (Bowman et al., 2015).
- **HANS Dataset:** McCoy et al. (2019) developed a dataset with the express purpose of scrutinizing the performance of models when confronted with sentences characterized by several types of distracting signals. These signals encompass the presence of lexical overlap, sub-sequences, and constituent heuristics between the corresponding hypotheses and premises.
- **Datasets Sensitive to Compositionality (ISCS):** As proposed by Nie et al. (2019), a softmax regression model was employed to utilize lexical features present in the premise and hypothesis sentences, thereby generating instances of misclassification. Here, the *Lexically Misleading Score* (LMS) denotes the predicted probability of the misclassified label. Adapting the approach of Liu et al. (2020b), we concentrated on the subsets possessing LMS values exceeding 0.7.
- **Not another Negation (NaN) NLI Dataset:** NaN dataset is developed to probe the capabilities of NLP models in comprehending sub-clausal negation (Truong et al., 2022).
- **Stress Test Datasets (ST):** Our analysis also incorporates various stress tests described by Naik et al. (2018) such as “word overlap” (ST-WO), “negation” (ST-NE), “length mismatch” (ST-LM), and “spelling errors” (ST-SE). Specifically, ST-WO aims to identify lexical overlap heuristics between the premise and hypothesis, ST-NE seeks to detect intense negative lexical cues in partial-input sentences, ST-LM aspires to create misleading predictions by artificially lengthening the premise using nonsensical phrases, and ST-SE employs spelling errors as a means to deceive the model.
- **Datasets Detected by Classifier (PICD):** In the approach proposed by Gururangan et al. (2018),

fastText was applied to hypothesis-only inputs. Subsequent instances from the SNLI test sets (Bowman et al., 2015) that could not be accurately classified were designated as ‘hard’ instances.

- **Surface Pattern Datasets (PISP):** Liu et al. (2020a) identified surface patterns that exhibit strong correlation with specific labels, thereby proposing adversarial test sets counteracting the implications of surface patterns. As suggested by Liu et al. (2020b), we employed their ‘hard’ instances extracted from the MultiNLI mismatched development set (Williams et al., 2018) as adversarial datasets.
- **Adversarial NLI (ANLI):** ANLI dataset (Nie et al., 2020) is a challenging resource created for training and testing models on NLI, featuring adversarial examples intentionally curated to obfuscate or mislead benchmark models, thereby increasing its challenge factor. This dataset is constructed in multiple rounds, with each subsequent round featuring human-created examples specifically designed to outsmart models trained on the previous rounds. In total, the dataset comprises three distinct rounds, specifically ANLI R1, ANLI R2, and ANLI R3, highlighting the layered complexity of this resource.
- **Quora Question Pairs (QQP):** QQP dataset (Wang et al., 2018) comprises pairs of questions sourced from the Quora community question-answering platform. The primary objective is to ascertain whether each question pair exhibits semantic equivalence.
- **Paraphrase Adversaries from Word Scrambling (PAWS):** The PAWS-QQP dataset (Zhang et al., 2019), derived from the QQP datasets, targets the intricate task of paraphrasing identification, emphasizing the differentiation of sentences that, despite high lexical similarity, convey distinct meanings. It incorporates adversarial examples generated via word scrambling, presenting a stringent assessment for NLP models.

B Meta-prompts for Generating Synthetic NLEs

Table 5 and 6 present the meta-prompts employed for producing NLEs utilizing ChatGPT in zero- and few-shot scenarios.

Meta-prompt for zero-shot generation
Assume that you’re an expert working on natural language inference tasks. Given a premise, a hypothesis, and the corresponding label. Please write a concise and precise reason to explain why the label is assigned to the example:
Meta-prompt for few-shot generation
Assume that you’re an expert working on natural language inference tasks. Given a premise, a hypothesis and the corresponding label. Please write a concise and precise reason to explain why the label is assigned to the example by following the provided examples:

Table 5: Meta-prompts used to generate NLEs via ChatGPT in zero- and few-shot scenarios for natural language inference tasks.

Meta-prompt for zero-shot generation
Assume that you’re an expert working on paraphrasing identification tasks. Given two sentences and the corresponding label. Please write a concise and precise reason to explain why the label is assigned to the example:
Meta-prompt for few-shot generation
Assume that you’re an expert working on paraphrasing identification tasks. Given two sentences and the corresponding label. Please write a concise and precise reason to explain why the label is assigned to the example by following the provided examples:

Table 6: Meta-prompts used to generate NLEs via ChatGPT in zero- and few-shot scenarios for paraphrasing identification tasks.

C Supplementary Studies

Using NLEs Generated by Vicuna and LLaMA2. Our research demonstrates that the integration of NLEs generated by ChatGPT significantly enhances the performance of X-ICL for five advanced LLMs. To assess the efficacy of these ChatGPT-generated NLEs, we explore the generation of synthetic NLEs using Vicuna and LLaMA2, ranked as the third and second-best models respectively. Likewise, these NLEs are generated in a few-shot setting, referred to herein as Vicuna_{few} and

Tasks	NLEs		
	fs-Vicuna	fs-LLaMA2	fs-ChatGPT
SNLI	62.9 (-5.0)	64.1 (-3.7)	65.0 (-2.9)
HANS	55.5 (-7.4)	67.4 (+4.5)	74.5 (+11.6)
ISCS	65.1 (+4.2)	63.6 (+2.7)	65.5 (+4.6)
NaN	62.6 (-1.6)	65.1 (+0.9)	66.3 (+2.1)
ST	59.5 (+2.2)	61.9 (+4.6)	64.8 (+7.5)
PICD	60.2 (-3.5)	60.8 (-2.9)	61.6 (-2.1)
PISP	66.0 (+11.0)	66.1 (+11.1)	66.0 (+11.0)
ANLI (R1)	66.1 (+9.1)	65.8 (+8.8)	64.9 (+7.9)
ANLI (R2)	55.4 (+6.5)	55.9 (+7.0)	55.5 (+6.6)
ANLI (R3)	49.6 (+10.8)	50.7 (+11.9)	52.0 (+13.2)
Average	60.3 (+3.8)	62.1 (+5.6)	63.5 (+6.9)

Table 7: ICL performance of Vicuna using (1) standard ICL without NLEs, (2) X-ICL with Vicuna-generated NLEs in a few-shot scenario: fs-Vicuna, (3) X-ICL with LLaMA2-generated NLEs in a few-shot scenario: fs-LLaMA2, (4) X-ICL with ChatGPT-generated NLEs in a few-shot scenario: fs-ChatGPT. Numbers in the parentheses represent differences compared to X-ICL (Human).

LLaMA2_{few}, respectively. To ensure a fair comparison, we employ Vicuna as the underlying model to evaluate fs-X-ICL(Vicuna), fs-X-ICL (LLaMA2), and fs-X-ICL (ChatGPT) on all studied datasets.

Our results, detailed in Table 7, highlight that X-ICL generally gains greater benefit from LLM-generated NLEs as opposed to those produced by humans. Meanwhile, fs-X-ICL (ChatGPT) consistently outperforms fs-X-ICL(Vicuna) and fs-X-ICL (LLaMA2) considerably, except for ANLI R1 and R2. These findings suggest that to harness the potential of AI-generated NLEs fully, the employment of a powerful LLM is integral.

Does model size matter? We have shown the efficacy of X-ICL across a range of LLMs of varying sizes. However, the variability in data and training processes among these models renders the applicability of our approach to smaller-scale models inconclusive, especially since the smaller models often exhibit less benefit from NLEs compared to larger models within the same family (Wei et al., 2022a). Therefore, we have evaluated our approach using three distinct sizes of LLaMA2 models: 7B, 13B, and 70B parameters.

Referring to Figure 6, one can find the performance of both ICL and X-ICL generally improves in correspondence with the escalation of model size, except for zs-X-ICL (ChatGPT). Moreover, the gap in performance between ICL and fs-X-ICL (ChatGPT) widens, indicating that models with greater capabilities derive increased benefits from

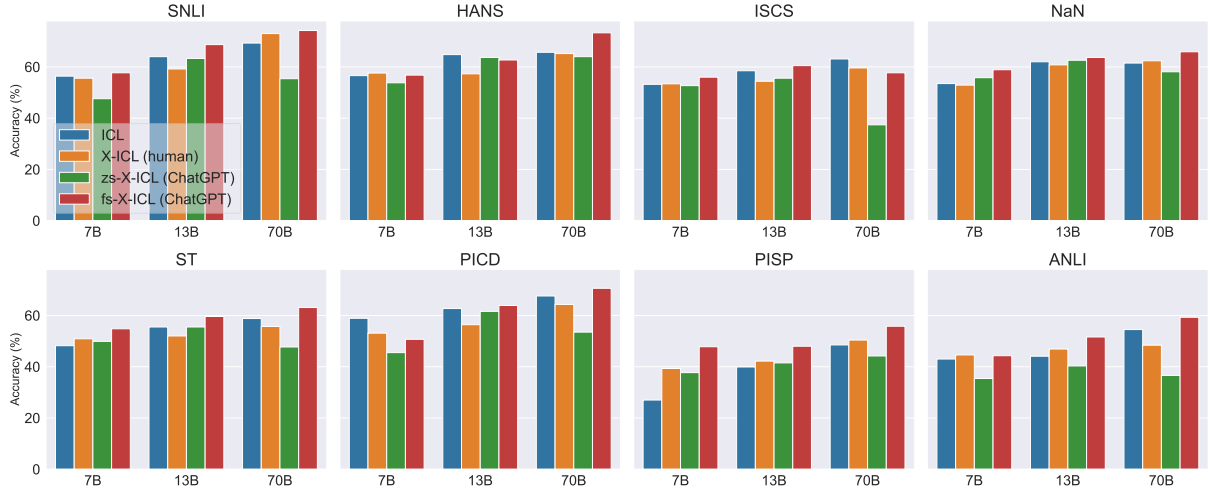


Figure 6: ICL performance of LLaMA2 (7B, 13B, 70B) using (1) standard ICL without NLEs, (2) X-ICL with human-written NLEs: X-ICL (Human), (3) X-ICL with ChatGPT-generated NLEs in a zero-shot scenario: zs-X-ICL (ChatGPT), (4) X-ICL with ChatGPT-generated NLEs in a few-shot scenario: fs-X-ICL (ChatGPT). ANLI is the average of R1, R2 and R3.

	NaN			PICD			ANLI (R1)			ANLI (R2)		
	e-SNLI	ANLI	$ \Delta $	e-SNLI	ANLI	$ \Delta $	e-SNLI	ANLI	$ \Delta $	e-SNLI	ANLI	$ \Delta $
ICL	70.0	69.4	0.6	64.0	64.1	0.1	52.6	62.4	9.7	43.9	51.7	7.8
fs-X-ICL (ChatGPT)	73.1	71.8	1.2	76.9	76.1	0.8	65.0	68.5	3.5	53.2	54.4	1.2

Table 8: Performance of ICL and fs-X-ICL (ChatGPT) employing e-SNLI and ANLI as prompts for testing NaN, PICD, and ANLI (R1/R2). $|\Delta|$ signifies the absolute difference in the performance outcomes when utilizing e-SNLI in contrast to ANLI. The backbone model is GPT3.5-turbo.

NLEs. This observation aligns with the results reported by Wei et al. (2022a).

Distribution Shift Prompting. Previous works indicate that X-ICL can potentially encourage LLMs to engage in deliberate thinking, a predominant factor responsible for substantial performance improvements over the standard ICL in complex reasoning tasks (Wei et al., 2022b). In addition, our findings have demonstrated a dramatic enhancement in the robustness of LLMs due to X-ICL, which contributes to significant improvements in ICL when applied to various adversarial datasets.

Moreover, a previous study established that upon understanding the concept underlying particular tasks, humans can address similar tasks despite a distribution shift (Scott, 1962). To explore the robustness of ICL and X-ICL against distribution shifts, we employ the e-SNLI dataset as the demonstration set for ANLI (R1/R2), while utilizing the ANLI training set for testing NaN and PICD. Due to its outstanding performance, we use GPT3.5-turbo as the backbone model.

As suggested in Table 8, for NaN and PICD, us-

ing e-SNLI as the prompt proves to be more effective than ANLI for both ICL and fs-X-ICL (ChatGPT). This improvement can be attributed to the distribution shift. Likewise, the distribution shift results in a noticeable distinction between e-SNLI and ANLI for ICL on ANLI (R1/R2). Nonetheless, incorporating NLEs enables fs-X-ICL (ChatGPT) to substantially reduce this gap, from 9.7 to 3.5 for ANLI (R1), and from 7.8 to 1.2 for ANLI (R2). This finding indicates that X-ICL may improve the robustness of LLMs in the face of distribution shifts.

Analysis on memorization LLMs such as ChatGPT have occasionally replicated instances from renowned benchmark datasets, including MNLI and BoolQ (Sainz et al., 2023). This unintentional ‘contamination’ might contribute to misconceptions regarding the superior performance of LLMs on these widespread benchmarks due to data memorization.

Following Carlini et al. (2023), we merge the premise and hypothesis of each test instance into a single sentence, using the first part as the prefix. If

an LLM could perfectly replicate the second part, we labeled the instance as ‘*extractable*’. Evaluating all studied models, we observe that the proportion of extractable instances is under 0.001% across all datasets and backbone models, indicating that the superior performance of LLMs might not be ascribed to memorization.

D Qualitative Analysis on NLEs

D.1 Qualitative Analysis on NLEs for Demonstration Set

We first conducted a qualitative analysis of NLEs generated by ChatGPT under zero- and few-shot scenarios, using the demonstration set as a basis. Note that each instance in the demonstration set has three distinct NLEs: (1) the zero-shot NLE from ChatGPT, (2) the few-shot NLE from ChatGPT, and (3) the human-written NLE. From these three NLEs per instance, one was randomly selected, and both the instance and the chosen NLE were incorporated into the evaluation set.

Subsequently, this evaluation set was rated independently by four authors on a 5-point Likert scale to assess the quality of the NLEs. The scale ranges were 1 (extremely dissatisfied), 2 (dissatisfied), 3 (neutral), 4 (satisfied), and 5 (extremely satisfied). Finally, we calculated the average scores for both ChatGPT-generated and human-written NLEs for each evaluator.

D.2 Qualitative Analysis on NLEs for Inference Set

We also conducted a qualitative analysis of NLEs generated by fs-X-ICL (ChatGPT), utilizing GPT3.5-turbo as the foundational model. A total of 280 randomly sampled, correctly predicted examples from fs-X-ICL (ChatGPT) were distributed evenly among seven evaluators. These evaluators were tasked to assess the quality of the NLE for each assigned instance, based on the premise-hypothesis pair and its corresponding correctly predicted label.

The evaluators were required to rate the quality of the NLE using the aforementioned 5-point Likert scale. In case of dissatisfaction, they were asked to identify the reason from a list of predefined factors, including:

- **template:** The NLE simply restates the input and employs it as a justification.

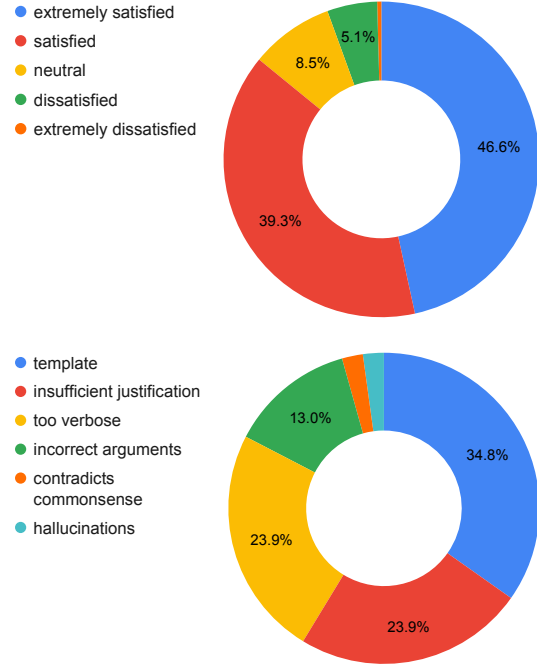


Figure 7: Human evaluation on ChatGPT-generated NLEs for the correct predictions from fs-X-ICL (ChatGPT). **Top:** distribution of satisfaction scores. **Bottom:** distribution of reasons for dissatisfaction.

- **insufficient justification:** The NLE requires more support for the prediction.
- **too verbose:** The NLE is overly detailed and includes unnecessary information.
- **incorrect arguments:** Despite the prediction being accurate, the NLE fails to support it due to erroneous arguments.
- **contradict commonsense:** The NLE is incorrect and contradicts commonsense.
- **hallucinations:** The NLE includes fabricated information.

According to Figure 7, 46.6% and 39.3% of NLEs are marked as ‘extremely satisfied’ and ‘satisfied’ respectively, constituting 85.9% of the total 280 NLE samples. This suggests a high-quality output from GPT3.5-turbo in general. As for the lower-quality NLEs, the primary reasons for dissatisfaction include ‘template’, ‘insufficient justification’, and ‘too verbose’. Interestingly, this suggests that, despite the expressed dissatisfaction, evaluators generally did not find incorrect justifications in most instances.

E Using AI Assitants

We use ChatGPT for the purpose of proofreading.

F Human-written NLEs for QQP

Given the absence of NLEs for the QQP dataset, we have conducted a randomized sampling of 20 instances from the QQP training set. For each selected instance, we crafted a corresponding NLE. The details of these QQP instances and their respective NLEs are presented in Table 9, 10, and 11.

<p>Q1: Is 6 foot 4 too tall as an ideal height for a man?</p> <p>Q2: My height is 5'6 and I'm 14 year old boy, my mom is 5'4 and my dad is 5'7. How tall will I be?</p> <p>Label: not duplicate</p> <p>NLE: Predicting future height given parents' heights concerns genetic factors of height, whereas ideal height for man concerns more about its social aspect.</p>
<p>Q1: Approximately how many hours have you spent on the internet till date?</p> <p>Q2: What amount of time do you spent on the Internet?</p> <p>Label: not duplicate</p> <p>NLE: Total number of hours spend on Internet till date not just depend on the average hours on internet per day, but also many other factors such as the age the user started using it.</p>
<p>Q1: What are the most ridiculous statements made by Donald Trump?</p> <p>Q2: My black friend supports Donald Trump, is that ridiculous?</p> <p>Label: not duplicate</p> <p>NLE: Asking the most ridiculous statement made by Donald Trump is different than asking why a supporter support him. A supporter can support him for other reasons.</p>
<p>Q1: "What is the origin of the phrase ""pipe dream""?"</p> <p>Q2: "How did the phrase ""toe head"" originate?"</p> <p>Label: not duplicate</p> <p>NLE: The two questions asked about the origin of two different words.</p>
<p>Q1: What is a good first programming language to learn?</p> <p>Q2: What is the most valuable programming language for the future to learn?</p> <p>Label: duplicate</p> <p>NLE: When picking a good first programming language to learn, people may consider the most valuable one language if they learn it for making money.</p>
<p>Q1: What is best way for earning money?</p> <p>Q2: How can I start making money? What are the best ways to make money?</p> <p>Label: duplicate</p> <p>NLE: Both questions ask about what are best ways to make money</p>
<p>Q1: Does the Indian education system need a reformation?</p> <p>Q2: Should the education system be changed in India? If so why or why not?</p> <p>Label: duplicate</p> <p>NLE: Both questions essentially inquire about the necessity and justification for changing the Indian education system.</p>
<p>Q1: What is the application of quantum physics?</p> <p>Q2: What are some applications of quantum physics?</p> <p>Label: duplicate</p> <p>NLE: The two questions both seek information about the practical use of quantum physics.</p>

Table 9: QQP instances and the corresponding NLEs.

<p>Q1: How is the word 'calumny' used in a sentence?</p> <p>Q2: How is the word 'mischievous' used in a sentence?</p> <p>Label: not duplicate</p> <p>NLE: The two questions ask about two different words with different meanings.</p>
<p>Q1: What are your views on the abolishment of 500 rupees note?</p> <p>Q2: How will the ban of Rs 500 and Rs 1000 notes affect Indian economy?</p> <p>Label: not duplicate</p> <p>NLE: The former question asks specifically about the abolishment of the Rs 500 note, while the latter asks about the Rs 500 and the Rs 1000 notes.</p>
<p>Q1: What are the valence electrons of titanium?</p> <p>Q2: What is the number of valence electrons in hydrogen? How is this determined?</p> <p>Label: not duplicate</p> <p>NLE: The former question asks about titanium, while the latter is about hydrogen.</p>
<p>Q1: Do movie actors get paid each time their movie is played on TV?</p> <p>Q2: Why are film actors so highly paid whereas scientists are paid relatively quite little?</p> <p>Label: not duplicate</p> <p>NLE: The former question asks some details about how actors get paid, while the latter asks about the gap between actor and scientist salaries.</p>
<p>Q1: How do I build an electromagnetic propulsion engine?</p> <p>Q2: How would I build a magnetic propulsion system?</p> <p>Label: duplicate</p> <p>NLE: Both question asks about building magnetic propulsion systems.</p>
<p>Q1: Why is salt water taffy candy imported in France?</p> <p>Q2: Why is Saltwater taffy candy imported in The Bahamas?</p> <p>Label: duplicate</p> <p>NLE: Both questions ask about the reasons behind importing salt water taffy candy.</p>
<p>Q1: Why do we call Java platform independent language when it still requires platform dependent JVM to get executed?</p> <p>Q2: How is the Java platform independent when we need to have JVM on every machine to run Java programs?</p> <p>Label: duplicate</p> <p>NLE: Both questions ask why do we call Java platform-independent, since it still depends on the availability of a JVM.</p>
<p>Q1: What are the various ways through which one can earn money online?</p> <p>Q2: How do you make easy money online?</p> <p>Label: duplicate</p> <p>NLE: Both questions ask how to make money online.</p>
<p>Q1: Does life get harder as you get older?</p> <p>Q2: Does life really get harder as you get older?</p> <p>Label: duplicate</p> <p>NLE: Both questions ask whether life does get harder as you get older.</p>

Table 10: QQP instances and the corresponding NLEs.

<p>Q1: Why can't some people think for themselves?</p> <p>Q2: Why don't people think for themselves?</p> <p>Label: not duplicate</p> <p>NLE: "some people" means not all people as the second question seems to imply</p>
<p>Q1: Why don't we use Solar Furnace to produce electricity?</p> <p>Q2: Why don't we make Solar Cars?</p> <p>Label: not duplicate</p> <p>NLE: using Solar Furnace you can produce some amount of electricity but it may not enough to power a whole car</p>
<p>Q1: What is an intuitive explanation of the fractional quantum Hall effect?</p> <p>Q2: What is an intuitive explanation of the Quantum Hall effect?</p> <p>Label: not duplicate</p> <p>NLE: fractional quantum Hall effect is different than the Quantum Hall effect, which refers to the integer quantum Hall effect</p>
<p>Q1: Can INTPs become successful entrepreneurs?</p> <p>Q2: I am business associate in tcs?</p> <p>Label: not duplicate</p> <p>NLE: completely different questions</p>
<p>Q1: How can I be like Sheldon Cooper?</p> <p>Q2: How do I become like Sheldon Cooper?</p> <p>Label: duplicate</p> <p>NLE: "be like" and "become like" someone is the same thing</p>
<p>Q1: What do people think about Anonymous?</p> <p>Q2: What do you think about the 'Anonymous' option on Quora?</p> <p>Label: duplicate</p> <p>NLE: "what do people think" and "what do you think" are usually used interchangeably</p>
<p>Q1: What's the meaning of life?</p> <p>Q2: "What is the meaning of ""Life""?"</p> <p>Label: duplicate</p> <p>NLE: same question with minor different spellings</p>
<p>Q1: What is it in for the Ibibo group employees with the Makemytrip merger / Buyout?</p> <p>Q2: How do Ibibo employees feel about MakeMyTrip acquiring Ibibo?</p> <p>Label: duplicate</p> <p>NLE: "the Makemytrip merger / Buyout" refers to "MakeMyTrip acquiring Ibibo" and "what is it in for the employees" means "how do the employees feel about"</p>

Table 11: QQP instances and the corresponding NLEs.