

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

Anonymous ACL submission

## Abstract

Recent studies have demonstrated that large language models (LLMs) excel in diverse tasks through in-context learning (ICL) facilitated by task-specific prompts and examples. However, the existing literature shows that ICL encounters performance deterioration when exposed to adversarial inputs. Enhanced performance has been observed when ICL is augmented with natural language explanations (NLEs) (we refer to it as X-ICL). Thus, this work investigates whether X-ICL can improve the robustness of LLMs on a suite of adversarial and challenging datasets covering natural language inference and paraphrasing identification. Moreover, we introduce a new approach to X-ICL by prompting an LLM (ChatGPT in our case) with few human-generated NLEs to produce further NLEs (we call it ChatGPT few-shot), which we show superior to both ChatGPT zero-shot and human-generated NLEs alone. We evaluate five popular LLMs (GPT3.5-turbo, LLaMa2, Vicuna, Zephyr, Mistral) and show that X-ICL with ChatGPT few-shot yields over 6% improvement over ICL. Furthermore, while prompt selection strategies were previously shown to improve ICL on in-distribution test sets significantly, we show that these strategies do not match the efficacy of the X-ICL paradigm in robustness-oriented evaluations.

## 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 unseen test data after observing a small number of demonstrations. Remarkably, they achieve this without the necessity for further training or any modifications to their underlying parameters. This novel learning paradigm is referred to as *in-context learning* (ICL, Brown et al., 2020; Rae et al., 2021). However, it has been noted that ICL struggles with the execution of complex tasks, such as arithmetic, commonsense, and

symbolic reasoning (Rae et al., 2021). To enhance the capability of ICL in solving tasks requiring complex reasoning, Wei et al. (2022b) draw inspiration from the extensive body of literature on natural language explanations (NLEs) to introduce a method denoted as the Chain-of-Thought (CoT) prompting. This method empowers LLMs to utilize human-written NLEs as a mechanism for deliberate thinking before delivering a prediction. Note that CoT and NLEs are interchangeable, describing the same concept. Since NLEs were introduced before CoT (Camburu et al., 2018; Hendricks et al., 2018), we used the former term in this paper. We denote the ICL equipped with NLEs as *X-ICL*. Albeit 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).

Similar to supervised learning, ICL demonstrates vulnerability to adversarial and misleading examples, causing a decline in performance (Wang et al., 2023a). Given that X-ICL promotes deliberate thinking in LLMs, we hypothesize that incorporating NLEs could enhance the resilience of LLMs against adversarial inputs, *aka* robustness. To this end, we leverage eight adversarial datasets to evaluate the added benefit of X-ICL to the robustness of LLMs.

Moreover, the effectiveness of X-ICL so far relies on human-written NLEs (Wei et al., 2022b), which usually require domain-specific expertise, thereby imposing constraints on its scalability. However, the advent of ChatGPT<sup>1</sup> uncovered a range of possibilities where LLMs can assist human annotators (Bang et al., 2023; Guo et al., 2023). Motivated by this development, we leverage ChatGPT (specifically, GPT3.5-turbo) to generate NLEs for examples from human-written NLEs. Following this generation step, four authors assess the quality of the human-written and ChatGPT-

<sup>1</sup><https://openai.com/blog/chatgpt>

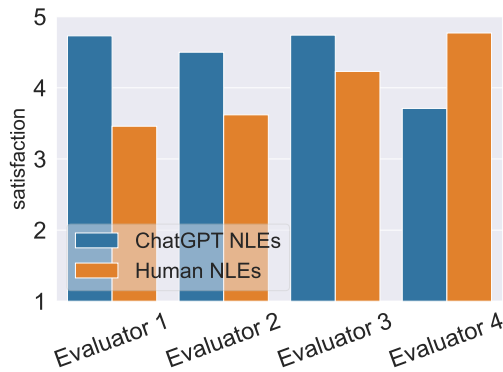


Figure 1: Human evaluation on ChatGPT-generated NLEs (ChatGPT NLEs) and human-written NLEs (Human NLEs). The satisfaction scores span from 1 (extremely dissatisfied) to 5 (extremely satisfied).

generated NLEs. As demonstrated in Figure 1, most of the evaluators (3 out of 4) exhibited a preference for the NLEs produced by ChatGPT over those crafted by humans. The details of this evaluation are presented in Appendix D.1.

In this paper, we evaluate the improvement in the robustness of LLMs provided by X-ICL in three regimes: utilizing NLEs generated by ChatGPT (generated in zero-shot and few-shot settings) and human-written 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 the challenging datasets.

Our experimental results suggest that X-ICL generally produces more accurate results than ICL on eight adversarial and challenging datasets. Furthermore, using few-shot ChatGPT-generated NLEs leads to more than 6% gains over ICL for the majority of the LLMs and datasets. The findings from our comprehensive study suggest that an integrated approach, combining human input with the capabilities of ChatGPT (*i.e.*, ChatGPT few-shot regime), provides a more effective solution than utilizing either human or ChatGPT (zero-shot) NLEs in isolation. Finally, while prompt-selection strategies (Gupta et al., 2023; Levy et al., 2023; Ye et al., 2023) considerably enhance ICL performance on in-distribution test sets, they are less effective on the adversarial datasets compared to the X-ICL approaches.

## 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 free-form natural language explanations (*i.e.*, NLEs) (Camburu et al., 2018; Narang et al., 2020; Wiegrefe and Marasovic, 2021). Our work concentrates on the latter category, namely, the generation of NLEs for justifying model predictions. Rajani et al. (2019) propose a two-stage training to improve the prediction performance for commonsense reasoning tasks. In their work, the first stage revolves around creating an NLE, which aids in label prediction during 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), medical (Kayser et al., 2022), and 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 continues to encounter difficulties with tasks that require reasoning abilities, including arithmetic, logical, and commonsense reasoning tasks (Rae et al., 2021; Srivastava et al., 2022). To augment the reasoning capabilities of LLMs, Wei et al. (2022b) introduced a method known as CoT prompting. This technique prompts an 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 an NLE/CoT before the prediction. Subsequently, Zhou et al. (2023) demonstrate that dividing complex problems into simpler sub-problems and addressing them sequentially improves the performance of CoT prompting. Additionally, Wang et al. (2023b) propose an alternative method that enhances CoT prompting by combining multiple diverse reasoning paths generated by LLMs, surpassing the performance of a greedy CoT prompting approach. However, these aforementioned methods need human-written NLEs as CoT. Instead, our ChatGPT zero-shot regime harnesses the power of an LLM to synthesize NLEs without

169 the need for human-written NLEs.

170 **Learning Robust Models.** Several works show  
171 that NLP models are prone to performance degra-  
172 dation when presented with adversarial datasets, a  
173 consequence of inherent artifacts or biases within  
174 the annotation of the training dataset (Naik et al.,  
175 2018; McCoy et al., 2019; Nie et al., 2020; Liu  
176 et al., 2020b). To mitigate biases within NLP mod-  
177 els, various strategies have been proposed, *e.g.*,  
178 initially training a weak model to recognize superfi-  
179 cial features, subsequently enforcing a target model  
180 to learn more robust and generalizable characteris-  
181 tics (He et al., 2019; Clark et al., 2019; Karimi Ma-  
182 habadi et al., 2020). Additionally, data augmenta-  
183 tion presents another viable option (Minervini and  
184 Riedel, 2018; Wu et al., 2021, 2022). Moreover,  
185 studies have shown that incorporation of rational-  
186 ization methodologies into supervised models can  
187 significantly enhance the models’ resilience against  
188 adversarial datasets (Chen et al., 2022; Stacey et al.,  
189 2022). Deviating from the precedent research, our  
190 study probes the robustness of X-ICL on eight ex-  
191 isting adversarial datasets.

### 192 3 Methodology

193 This section first outlines the workflow of X-ICL.  
194 Subsequently, the focus shifts to detailing how an  
195 LLM, specifically ChatGPT, can be used to gener-  
196 ate an NLE for a labeled instance.

#### 197 3.1 ICL with NLEs (X-ICL)

198 LLMs greatly enhance their performance across  
199 various reasoning tasks when supplied with human-  
200 written NLEs (Wei et al., 2022b,a). We can define  
201 X-ICL as follows:

$$202 \arg \max_{(\mathbf{r}', \mathbf{y}') \in \mathbb{R} \times \mathbb{Y}} P_{\theta} \left( (\mathbf{r}', \mathbf{y}') | (\mathbf{x}_i, \mathbf{r}_i, \mathbf{y}_i)_{i=1}^k, \dots, (\mathbf{x}') \right),$$

203 where  $\mathbf{x}'$  denotes an unlabeled instance,  $\mathbf{r} \in \mathbb{R}$  rep-  
204 represents the corresponding NLE, and  $\mathbf{y} \in \mathbb{Y}$  denotes  
205 the target label, where  $\mathbb{R}$  is the space of all NLEs  
206 and  $\mathbb{Y}$  is the set of possible labels for a given dataset.  
207 The objective of X-ICL is to maximize the likeli-  
208 hood of generating the optimal NLE,  $\mathbf{r}' \in \mathbb{R}$ , and  
209 its corresponding label,  $\mathbf{y}' \in \mathbb{Y}$ , given a demonstra-  
210 tion set  $(\mathbf{x}_i, \mathbf{r}_i, \mathbf{y}_i)_{i=1}^k$  and an unlabeled instance  $\mathbf{x}'$ .  
211 Consequently, this prompts the LLM to produce  
212 the most plausible NLE and label combination.

#### 213 3.2 Generating NLEs via ChatGPT

214 In existing X-ICL works, human-written NLEs  $\mathbf{r}$   
215 were used for the instances within the demonstra-

216 tion set. Instead, in this work, we opt for the NLEs  
217 synthesized via ChatGPT (or GPT3.5-turbo). This  
218 preference is driven by noting that NLEs produced  
219 by ChatGPT tend to receive higher approval ratings  
220 from human evaluators, as indicated in Figure 1.  
221 We argue that this preference will boost the per-  
222 formance of X-ICL. The methods utilized for the  
223 generation of NLEs are outlined below.

224 **Few-shot prompting for NLEs** Our methodol-  
225 ogy, also shown in Figure 2, initiates by leveraging  
226 a set of labeled instances, each accompanied by  
227 a human-crafted NLE, to prompt ChatGPT. The  
228 primary aim is to encourage the LLMs to generate  
229 a correct NLE (*i.e.*, ground-truth arguments) for  
230 the correctly-predicted answer for a test instance.  
231 The NLE is generated as follows:

$$232 \arg \max_{\mathbf{r}' \in \mathbb{R}} P_{\theta}(\mathbf{r}' | \mathbf{s}, (\mathbf{x}_j, \mathbf{y}_j, \mathbf{r}_j)_{j=1}^m, \dots, (\mathbf{x}', \mathbf{y}')),$$

233 where  $\mathbf{s}$  is a meta-prompt representing the task.  
234 For the details of the meta-prompt, please refer to  
235 Appendix B.

236 **Zero-shot prompting for NLEs** We further ex-  
237 tend our approach to situations where human-  
238 written NLEs are absent, which is generally more  
239 prevalent across most datasets. In this context,  
240 ChatGPT is prompted to generate an NLE for a  
241 labeled instance devoid of any pre-existing exam-  
242 ples with NLEs. The objective bears a resemblance  
243 to Equation (1), albeit without the inclusion of the  
244 demonstration set  $(\mathbf{x}_j, \mathbf{y}_j, \mathbf{r}_j)_{j=1}^m$ .

245 Notably, the NLEs generated by the aforemen-  
246 tioned approaches can be seamlessly integrated  
247 into the existing X-ICL framework as delineated  
248 in Section 3.1. These are referred to as X-ICL  
249 (ChatGPT<sub>few</sub>) and X-ICL (ChatGPT<sub>zero</sub>), respec-  
250 tively.<sup>2</sup>

## 251 4 Experiments

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

### 254 4.1 Experimental Setup

255 **Tasks and datasets** We consider the Natural Lan-  
256 guage Inference (NLI) and paraphrasing identifi-  
257 cation tasks as our testbed. To ascertain the ro-  
258 bustness of LLMs when employing the proposed

<sup>2</sup>In addition, we explore the application of two other  
widely-used, open-source LLMs for the generation of NLEs.  
Detailed results of these experiments are provided in Ap-  
pendix C.

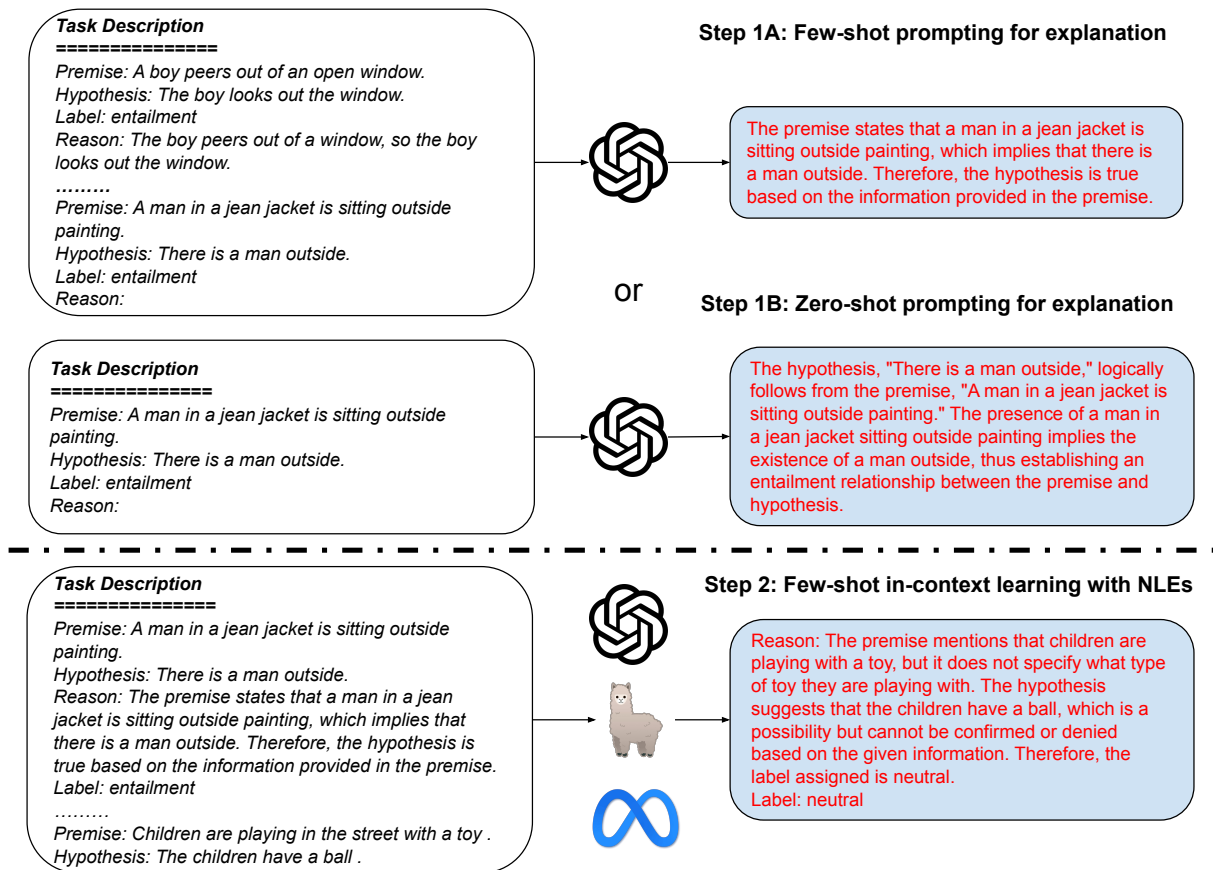


Figure 2: Our approach of using ChatGPT-generated NLEs for ICL consists of two steps: (1) prompt an LLM in a few-shot or zero-shot manner to generate NLEs for new instances; (2) prompt LLMs using ICL with the NLEs generated in step 1.

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 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 (Jiang et al., 2023), (2) Zephyr (Tunstall et al., 2023), (3) Vicuna (Chiang et al., 2023), (4) LLaMA2 (Touvron et al., 2023), and (5) GPT3.5-turbo. Specifically, Mistral and Zephyr are models with 7B parameters. For Vicuna and LLaMA2, we use 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, barring 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 (Camburu et al., 2018) is a modified version of the SNLI dataset, enhanced with NLEs written by humans. In the case of the QQP and PAWS datasets, the QQP dataset is utilized as the demonstration set, including NLEs contributed by four authors.

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.



Models	Methods	SNLI	HANS	ISCS	NaN	ST	PICD	PISP	ANLI	QQP	PAWS	Avg.
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 <sup>▽</sup> ±2.5	58.6 <sup>▽</sup> ±2.9	51.7 <sup>▼</sup> ±4.0	56.9 ±3.3	35.8 <sup>▼</sup> ±6.7	43.9 <sup>▼</sup> ±1.7	69.9 ±0.8	66.4 ±1.5	53.5
	X-ICL (ChatGPT <sub>zero</sub> )	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
	X-ICL (ChatGPT <sub>few</sub> )	<b>61.8</b> ±3.1	<b>58.2<sup>▼</sup></b> ±2.5	<b>57.2<sup>▼</sup></b> ±2.2	<b>62.4<sup>▼</sup></b> ±2.6	<b>55.2<sup>▼</sup></b> ±1.5	<b>59.2</b> ±2.7	<b>47.6<sup>▼</sup></b> ±1.8	<b>46.9<sup>▼</sup></b> ±2.3	<b>70.3</b> ±1.1	<b>72.5<sup>▽</sup></b> ±1.3	<b>57.1</b>
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 <sup>▼</sup> ±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 <sup>▽</sup> ±2.3	59.3
	X-ICL (ChatGPT <sub>zero</sub> )	67.2 ±3.9	72.7 ±2.6	60.4 ±3.3	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
	X-ICL (ChatGPT <sub>few</sub> )	<b>74.2<sup>▼</sup></b> ±3.6	<b>77.4<sup>▼</sup></b> ±2.2	<b>67.0</b> ±1.6	<b>67.7</b> ±2.3	<b>69.3<sup>▼</sup></b> ±1.5	<b>70.0<sup>▼</sup></b> ±2.1	<b>65.6<sup>▼</sup></b> ±2.5	<b>52.1<sup>▽</sup></b> ±2.8	<b>77.3</b> ±0.9	<b>61.5<sup>▽</sup></b> ±1.0	<b>65.5</b>
Vicuna 30B	ICL	65.2 ±2.7	69.4 ±1.2	62.7 ±0.9	61.4 ±3.5	58.7 ±0.8	<b>67.1</b> ±1.6	50.9 ±1.3	50.0 ±2.6	<b>81.8</b> ±0.5	69.7 ±2.6	61.4
	X-ICL (Human)	<b>67.8</b> ±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
	X-ICL (ChatGPT <sub>zero</sub> )	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
	X-ICL (ChatGPT <sub>few</sub> )	65.0 ±3.1	<b>74.5<sup>▽</sup></b> ±4.4	<b>65.5<sup>▽</sup></b> ±1.6	<b>66.3<sup>▽</sup></b> ±1.1	<b>64.8<sup>▼</sup></b> ±1.8	61.6 ±8.9	<b>65.9<sup>▼</sup></b> ±4.7	<b>57.5<sup>▼</sup></b> ±1.3	78.6 ±1.7	<b>70.0</b> ±3.3	<b>65.4</b>
LLaMA2 70B	ICL	69.3 ±1.2	65.7 ±3.4	<b>63.1</b> ±1.6	61.5 ±2.3	58.8 ±4.4	67.6 ±3.0	48.5 ±7.3	54.2 ±2.9	<b>80.8</b> ±0.6	44.5 ±2.9	60.3
	X-ICL (Human)	73.0 <sup>▼</sup> ±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
	X-ICL (ChatGPT <sub>zero</sub> )	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
	X-ICL (ChatGPT <sub>few</sub> )	<b>74.2<sup>▼</sup></b> ±2.5	<b>73.3<sup>▼</sup></b> ±8.5	57.7 ±1.2	<b>65.9<sup>▽</sup></b> ±3.2	<b>63.1<sup>▽</sup></b> ±3.7	<b>70.6<sup>▽</sup></b> ±6.5	<b>55.8<sup>▼</sup></b> ±5.9	<b>59.2<sup>▼</sup></b> ±1.6	77.6 ±0.6	<b>46.5<sup>▽</sup></b> ±1.9	<b>63.6</b>
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	<b>81.5</b> ±0.3	42.9 ±2.8	62.4
	X-ICL (Human)	<b>78.0<sup>▼</sup></b> ±1.7	71.0 ±1.7	69.0 <sup>▽</sup> ±1.2	70.5 ±2.2	65.7 <sup>▽</sup> ±2.2	72.7 <sup>▼</sup> ±1.3	59.3 <sup>▽</sup> ±1.9	59.8 <sup>▽</sup> ±2.3	76.0 ±3.9	53.4 <sup>▼</sup> ±5.3	66.2
	X-ICL (ChatGPT <sub>zero</sub> )	71.9 ±2.7	71.6 ±0.8	68.4 <sup>▽</sup> ±0.3	70.2 ±0.0	67.6 <sup>▽</sup> ±1.3	67.7 <sup>▽</sup> ±4.1	61.7 <sup>▼</sup> ±1.9	<b>60.4<sup>▼</sup></b> ±2.0	80.4 ±0.8	51.2 <sup>▼</sup> ±3.1	66.0
	X-ICL (ChatGPT <sub>few</sub> )	75.5 <sup>▽</sup> ±2.8	<b>76.0<sup>▼</sup></b> ±2.0	<b>74.9<sup>▼</sup></b> ±0.1	<b>73.1<sup>▼</sup></b> ±1.4	<b>73.3<sup>▼</sup></b> ±0.4	<b>76.9<sup>▼</sup></b> ±0.4	<b>75.5<sup>▼</sup></b> ±3.0	59.6 <sup>▽</sup> ±1.8	79.0 ±1.7	<b>54.0<sup>▼</sup></b> ±2.6	<b>69.7</b>

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: X-ICL (ChatGPT<sub>zero</sub>), (4) X-ICL with ChatGPT-generated NLEs in a few-shot scenario: X-ICL (ChatGPT<sub>few</sub>). 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: ▼ (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.

**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.

Firstly, the results reveal a predictable outcome for both scenarios, namely with and without using X-ICL. As the models’ abilities escalate, an upward

trajectory can be discerned in the average accuracy. This progression is evident when comparing the least potent model, exemplified by Mistral, to the highest-performing 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.

In the context of X-ICL (ChatGPT<sub>few</sub>), the evi-

dence points to a commanding lead in all evaluated tasks on both the Mistral and Zephyr, surpassing the results of ICL and X-ICL (Human) by margins of at least 5.7% and 3.6%, respectively. Despite the notable improvement on ICL when employing GPT3.5-turbo in comparison to other LLMs, X-ICL (ChatGPT<sub>few</sub>) offers substantially additional gains, with an increase in absolute accuracy between 11%-24% on tasks such as ISCS, ST, PICD, PISP and PAWS. In essence, X-ICL augments LLM performance on the in-distribution test and bolsters LLMs’ robustness in the face of adversarial test sets.

Remarkably, despite the predominant preference of human evaluators for NLEs generated by ChatGPT over those written by humans, X-ICL (ChatGPT<sub>zero</sub>) 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 X-ICL (ChatGPT<sub>zero</sub>) is outperformed by ICL for 4 out of the 5 LLMs analyzed, especially on LLaMA2. This apparent discrepancy between human preferences and LLM performance strongly underlines the necessity for additional investigations to enhance our understanding of this intriguing phenomenon. Since this investigation deserves a thorough and systematic study, we leave it for future work.

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

### 4.3 Impacts of NLEs

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

**Data selection vs. X-ICL.** The efficacy of ICL in LLMs is significantly influenced by the demonstrations provided, as the model depends on these demonstrations to comprehend and address the test instances (Zhao et al., 2021; Liu et al., 2022; Lu et al., 2022). Consequently, a spectrum of research has been directed towards optimizing data selection techniques that curate ICL demonstrations from a pertinent pool of candidate data in relation to the test instances (Gupta et al., 2023; Levy et al.,

Models	Methods	SNLI	AdvNLI	$\Delta$
Mistral	ICL	59.8	45.1	14.7
	X-ICL (ChatGPT <sub>few</sub> )	61.8	<b>53.4</b>	<b>8.4</b>
	COSINE	67.9	46.0	21.9
	BM25	65.2	44.2	21.0
	SET-BSR	<b>77.6</b>	52.2	25.4
Zephyr	ICL	67.1	57.2	<b>9.9</b>
	X-ICL (ChatGPT <sub>few</sub> )	74.2	<b>63.7</b>	10.5
	COSINE	77.0	55.6	21.4
	BM25	70.1	53.7	16.4
	SET-BSR	<b>79.9</b>	59.7	20.2
Vicuna	ICL	65.2	57.8	7.4
	X-ICL (ChatGPT <sub>few</sub> )	65.0	<b>63.5</b>	<b>1.5</b>
	COSINE	72.5	53.6	18.9
	BM25	67.2	52.2	15.0
	SET-BSR	<b>79.5</b>	56.4	23.1
LLaMA2	ICL	69.3	58.7	<b>10.6</b>
	X-ICL (ChatGPT <sub>few</sub> )	74.2	<b>62.7</b>	11.5
	COSINE	71.9	57.3	14.6
	BM25	70.8	55.6	15.2
	SET-BSR	<b>76.7</b>	59.2	17.5
GPT3.5-turbo	ICL	71.9	61.4	10.5
	X-ICL (ChatGPT <sub>few</sub> )	75.5	<b>69.8</b>	<b>5.6</b>
	COSINE	75.0	58.1	16.9
	BM25	71.4	56.0	15.4
	SET-BSR	<b>77.4</b>	59.5	17.9

Table 2: Performance of ICL, X-ICL (ChatGPT<sub>few</sub>) and three data selection approaches on SNLI and AdvNLI (i.e., 7 adversarial test sets).  $\Delta$  indicates the difference between SNLI and adversarial NLI test sets. We report the average performance over all adversarial test sets.

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 examine the performance of X-ICL (ChatGPT<sub>few</sub>) in relation 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), coupled with strategies to ensure information coverage at the set level, to promote the choice of informative and diverse demonstration sets (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 process that is both computationally demanding and cost-inefficient for generating NLEs for the full dataset. Therefore, our analysis is confined

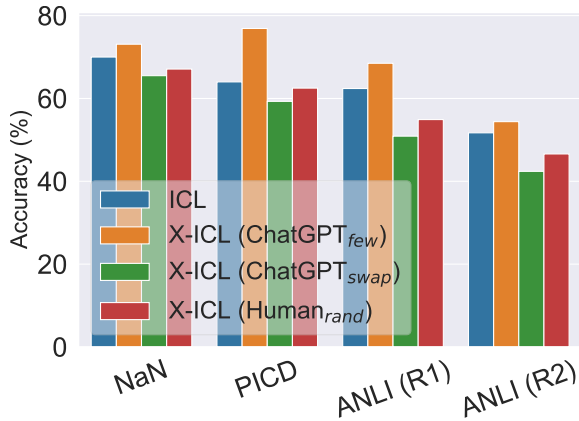


Figure 3: ICL performance of GPT3.5-turbo using (1) standard ICL without NLEs, (2) X-ICL with ChatGPT-generated NLEs in a few-shot scenario: X-ICL (ChatGPT<sub>few</sub>), (3) X-ICL with ChatGPT-generated NLEs, where the NLEs of the prompt are swapped and do not match the instances: X-ICL (ChatGPT<sub>swap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

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 X-ICL (ChatGPT<sub>few</sub>). 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, 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

X-ICL (ChatGPT<sub>swap</sub>). 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<sub>rand</sub>).

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 X-ICL (ChatGPT<sub>swap</sub>) when compared to X-ICL (ChatGPT<sub>few</sub>). This decline is discernible across various datasets, including NaN, PICD, and ANLI (R1/R2).<sup>3</sup> It is also shown that an irrelevant and arbitrary NLE triggers a performance reduction within the X-ICL framework. Furthermore, the efficiency of both X-ICL (ChatGPT<sub>swap</sub>) and X-ICL (Human<sub>rand</sub>) substantially lags behind that of ICL. Therefore, it can be inferred that the efficacy of the X-ICL (ChatGPT<sub>few</sub>) hinges on providing an accurate and relevant NLE.

#### 4.4 Supplementary Studies

**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 4, one can find the performance of both ICL and X-ICL generally improves in correspondence with the escalation of model size, except for X-ICL (ChatGPT<sub>zero</sub>). Moreover, the gap in performance between ICL and X-ICL (ChatGPT<sub>few</sub>) widens, indicating that models with greater capabilities derive increased benefits from 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,

<sup>3</sup>Similar patterns have been detected in other datasets

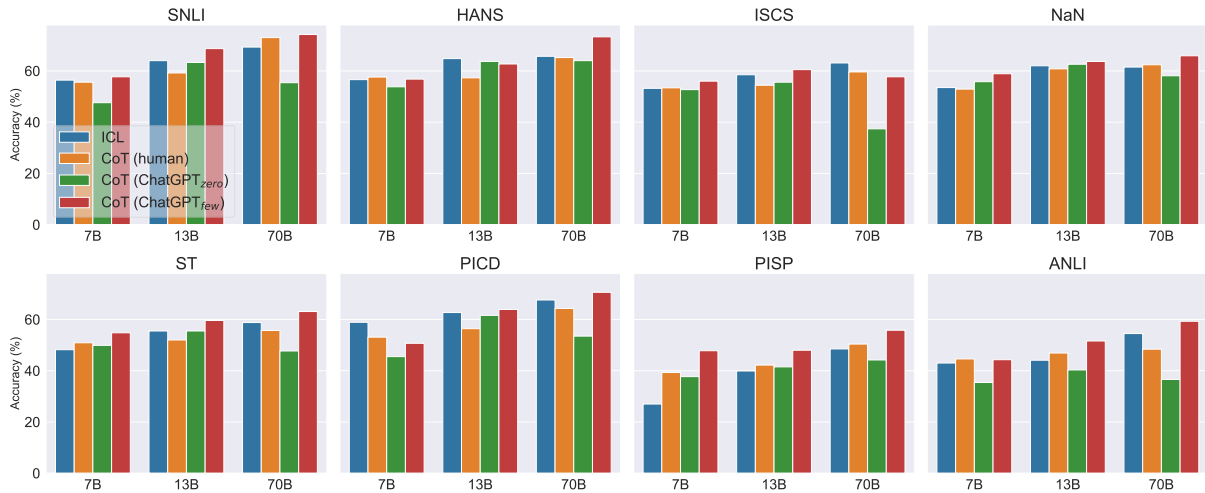


Figure 4: ICL performance of LLaMA2 (7B, 13B, 70B) using (1) standard ICL without NLEs, (2) X-ICL with human-written NLEs: X-ICL (Human), (3) X-ICL with ChatGPT-generated NLEs in a zero-shot scenario: X-ICL (ChatGPT<sub>zero</sub>), (4) X-ICL with ChatGPT-generated NLEs in a few-shot scenario: X-ICL (ChatGPT<sub>few</sub>). 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
X-ICL (ChatGPT <sub>few</sub> )	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 3: Performance of ICL and X-ICL (ChatGPT<sub>few</sub>) 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.

485 which contributes to significant improvements in  
486 ICL when applied to various adversarial datasets.

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

497 As suggested in Table 3, for NaN and PICD,  
498 using e-SNLI as the prompt proves to be more  
499 effective than ANLI for both ICL and X-ICL  
500 (ChatGPT<sub>few</sub>). This improvement can be at-  
501 tributed to the distribution shift. Likewise, the  
502 distribution shift results in a noticeable distinc-  
503 tion between e-SNLI and ANLI for ICL on ANLI  
504 (R1/R2). Nonetheless, incorporating NLEs enables  
505 X-ICL (ChatGPT<sub>few</sub>) to substantially reduce this  
506 gap, from 9.7 to 3.5 for ANLI (R1), and from 7.8  
507 to 1.2 for ANLI (R2). This finding indicates that  
508 X-ICL may improve the robustness of LLMs in the

face of distribution shifts.

## 5 Summary and Outlook

509 We introduced a simple yet effective method called  
510 X-ICL (ChatGPT<sub>few</sub>), leveraging human-written  
511 NLEs to generate synthetic NLEs by prompt-  
512 ing ChatGPT. X-ICL (ChatGPT<sub>few</sub>) significantly  
513 boosts accuracy across various adversarial datasets  
514 and five LLMs, compared to standard in-context  
515 learning and X-ICL using human-written NLEs.  
516 Additionally, our analysis revealed that data selec-  
517 tion methodologies may exhibit overfitting within  
518 the in-distribution dataset, thus potentially failing  
519 to extend to unseen or adversarial datasets. In con-  
520 trast, our approach employing NLEs has shown  
521 consistent performance in both in-distribution and  
522 adversarial contexts. Our work paves the way for  
523 more robust performance and enhanced explain-  
524 ability capabilities of LLMs.

## Limitations

525 One limitation of X-ICL might be the observed  
526 lack of fidelity in the NLEs generated by LLMs,  
527 despite their capability to provide accurate answers.  
528  
529  
530



531	These NLEs may sometimes include unfaithful or hallucinated information, which if relied upon by	Oana-Maria Camburu, Tim Rocktäschel, Thomas	584
532	users for model trust, can lead to severe implica-	Lukasiewicz, and Phil Blunsom. 2018. e-snli: Natu-	585
533	tions. Testing and enhancing the faithfulness of	ral language inference with natural language expla-	586
534	NLEs is a challenging open question (Atanasova	nations. <i>Advances in Neural Information Processing</i>	587
535	et al., 2023). In this work, we show that X-ICL im-	<i>Systems</i> , 31.	588
536	proves robustness, but we do not advocate for the	Nicholas Carlini, Daphne Ippolito, Matthew Jagielski,	589
537	usage of the generated NLEs as faithful explana-	Katherine Lee, Florian Tramèr, and Chiyuan Zhang.	590
538	tions without further testing. Second, our approach	2023. <a href="#">Quantifying memorization across neural lan-</a>	591
539	exhibited promising results when tested against ad-	<a href="#">guage models</a> . In <i>The Eleventh International Confer-</i>	592
540	versarial datasets in two notable NLP tasks: natural	<a href="#">ence on Learning Representations</a> .	593
541	language inference and paraphrasing identification.	Howard Chen, Jacqueline He, Karthik Narasimhan, and	594
542	However, further research is required to examine	Danqi Chen. 2022. Can rationalization improve ro-	595
543	the performance of LLMs and their generalizability	bustness? In <i>Proceedings of the 2022 Conference</i>	596
544	across diverse NLP tasks in the context of adversar-	<i>of the North American Chapter of the Association</i>	597
545	ial examples.	<i>for Computational Linguistics: Human Language</i>	598
546		<i>Technologies</i> , pages 3792–3805.	599
547	<b>References</b>	Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,	600
548	David Alvarez-Melis and Tommi Jaakkola. 2017. A	Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan	601
549	causal framework for explaining the predictions of	Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion	602
550	black-box sequence-to-sequence models. In <i>Proce-</i>	Stoica, and Eric P. Xing. 2023. <a href="#">Vicuna: An open-</a>	603
551	<i>edings of the 2017 Conference on Empirical Methods</i>	<a href="#">source chatbot impressing gpt-4 with 90%* chatgpt</a>	604
552	<i>in Natural Language Processing</i> , pages 412–421,	<a href="#">quality</a> .	605
553	Copenhagen, Denmark. Association for Computa-	Christopher Clark, Mark Yatskar, and Luke Zettlemoyer.	606
554	tional Linguistics.	2019. <a href="#">Don’t take the easy way out: Ensemble based</a>	607
555	Pepa Atanasova, Oana-Maria Camburu, Christina Li-	<a href="#">methods for avoiding known dataset biases</a> . In <i>Pro-</i>	608
556	oma, Thomas Lukasiewicz, Jakob Grue Simonsen,	<i>ceedings of the 2019 Conference on Empirical Meth-</i>	609
557	and Isabelle Augenstein. 2023. Faithfulness Tests for	<i>ods in Natural Language Processing and the 9th In-</i>	610
558	Natural Language Explanations. In <i>ACL</i> .	<i>ternational Joint Conference on Natural Language</i>	611
559	Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-	<i>Processing (EMNLP-IJCNLP)</i> , pages 4069–4082,	612
560	liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei	Hong Kong, China. Association for Computational	613
561	Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu,	Linguistics.	614
562	and Pascale Fung. 2023. A multitask, multilingual,	Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang,	615
563	multimodal evaluation of chatgpt on reasoning, hal-	Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng	616
564	lucination, and interactivity. <i>CoRR</i> , abs/2302.04023.	Wu. 2023. How close is chatgpt to human experts?	617
565	Samuel R. Bowman, Gabor Angeli, Christopher Potts,	comparison corpus, evaluation, and detection. <i>CoRR</i> ,	618
566	and Christopher D. Manning. 2015. A large an-	abs/2301.07597.	619
567	notated corpus for learning natural language infer-	Shivanshu Gupta, Sameer Singh, and Matt Gardner.	620
568	ence. In <i>EMNLP</i> , pages 632–642. The Association	2023. Coverage-based example selection for in-	621
569	for Computational Linguistics.	context learning. <i>arXiv preprint arXiv:2305.14907</i> .	622
570	Tom Brown, Benjamin Mann, Nick Ryder, Melanie	Suchin Gururangan, Swabha Swayamdipta, Omer Levy,	623
571	Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind	Roy Schwartz, Samuel Bowman, and Noah A. Smith.	624
572	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	2018. <a href="#">Annotation artifacts in natural language infer-</a>	625
573	Askell, Sandhini Agarwal, Ariel Herbert-Voss,	<a href="#">ence data</a> . In <i>Proceedings of the 2018 Conference of</i>	626
574	Gretchen Krueger, Tom Henighan, Rewon Child,	<i>the North American Chapter of the Association for</i>	627
575	Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens	<i>Computational Linguistics: Human Language Tech-</i>	628
576	Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-	<i>nologies, Volume 2 (Short Papers)</i> , pages 107–112,	629
577	teusz Litwin, Scott Gray, Benjamin Chess, Jack	New Orleans, Louisiana. Association for Computa-	630
578	Clark, Christopher Berner, Sam McCandlish, Alec	tional Linguistics.	631
579	Radford, Ilya Sutskever, and Dario Amodei. 2020.	Peter Hase, Shiyue Zhang, Harry Xie, and Mohit Bansal.	632
580	<a href="#">Language models are few-shot learners</a> . In <i>Ad-</i>	2020. <a href="#">Leakage-adjusted simulatability: Can models</a>	633
581	<i>vances in Neural Information Processing Systems</i> ,	<a href="#">generate non-trivial explanations of their behavior</a>	634
582	volume 33, pages 1877–1901. Curran Associates,	<a href="#">in natural language?</a> In <i>Findings of the Association</i>	635
583	Inc.	<i>for Computational Linguistics: EMNLP 2020</i> , pages	636
		4351–4367, Online. Association for Computational	637
		Linguistics.	638
		He He, Sheng Zha, and Haohan Wang. 2019. <a href="#">Unlearn</a>	639
		<a href="#">dataset bias in natural language inference by fitting</a>	640

641	<a href="#">the residual</a> . In <i>Proceedings of the 2nd Workshop on Deep Learning Approaches for Low-Resource NLP (DeepLo 2019)</i> , pages 132–142, Hong Kong, China. Association for Computational Linguistics.	
642		
643		
644		
645	Lisa Anne Hendricks, Ronghang Hu, Trevor Darrell, and Zeynep Akata. 2018. Grounding visual explanations. In <i>Proceedings of the European Conference on Computer Vision (ECCV)</i> .	
646		
647		
648		
649	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	
650		
651		
652		
653		
654	Rabeeh Karimi Mahabadi, Yonatan Belinkov, and James Henderson. 2020. <a href="#">End-to-end bias mitigation by modelling biases in corpora</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8706–8716, Online. Association for Computational Linguistics.	
655		
656		
657		
658		
659		
660	Maxime Kayser, Oana-Maria Camburu, Leonard Salewski, Cornelius Emde, Virginie Do, Zeynep Akata, and Thomas Lukasiewicz. 2021. e-ViL: A dataset and benchmark for natural language explanations in vision-language tasks. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 1244–1254.	
661		
662		
663		
664		
665		
666		
667	Maxime Kayser, Cornelius Emde, Oana-Maria Camburu, Guy Parsons, Bartłomiej Papież, and Thomas Lukasiewicz. 2022. Explaining chest x-ray pathologies in natural language. In <i>Medical Image Computing and Computer Assisted Intervention – MICCAI 2022</i> , pages 701–713, Cham. Springer Nature Switzerland.	
668		
669		
670		
671		
672		
673		
674	Jinkyu Kim, Anna Rohrbach, Trevor Darrell, John F. Canny, and Zeynep Akata. 2018. <a href="#">Textual explanations for self-driving vehicles</a> . <i>CoRR</i> , abs/1807.11546.	
675		
676		
677		
678	Itay Levy, Ben Bogin, and Jonathan Berant. 2023. <a href="#">Diverse demonstrations improve in-context compositional generalization</a> . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1401–1422, Toronto, Canada. Association for Computational Linguistics.	
679		
680		
681		
682		
683		
684		
685	Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen, Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian, Baolin Peng, Yi Mao, et al. 2022. Explanations from large language models make small reasoners better. <i>arXiv preprint arXiv:2210.06726</i> .	
686		
687		
688		
689		
690	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. <a href="#">What makes good in-context examples for GPT-3?</a> In <i>Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures</i> , pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.	
691		
692		
693		
694		
695		
696		
697		
	Tianyu Liu, Zheng Xin, Baobao Chang, and Zhifang Sui. 2020a. <a href="#">HypoNLI: Exploring the artificial patterns of hypothesis-only bias in natural language inference</a> . In <i>Proceedings of the Twelfth Language Resources and Evaluation Conference</i> , pages 6852–6860, Marseille, France. European Language Resources Association.	698 699 700 701 702 703 704
	Tianyu Liu, Zheng Xin, Xiaoan Ding, Baobao Chang, and Zhifang Sui. 2020b. <a href="#">An empirical study on model-agnostic debiasing strategies for robust natural language inference</a> . In <i>Proceedings of the 24th Conference on Computational Natural Language Learning</i> , pages 596–608, Online. Association for Computational Linguistics.	705 706 707 708 709 710 711
	Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. <a href="#">Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.	712 713 714 715 716 717 718 719
	Bodhisattwa Prasad Majumder, Oana-Maria Camburu, Thomas Lukasiewicz, and Julian McAuley. 2022. Knowledge-grounded self-rationalization via extractive and natural language explanations. In <i>Proceedings of the 39th International Conference on Machine Learning</i> , volume 162 of <i>Proceedings of Machine Learning Research</i> , pages 14786–14801. PMLR.	720 721 722 723 724 725 726
	Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. <a href="#">Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3428–3448, Florence, Italy. Association for Computational Linguistics.	727 728 729 730 731 732
	Pasquale Minervini and Sebastian Riedel. 2018. <a href="#">Adversarially regularising neural NLI models to integrate logical background knowledge</a> . In <i>Proceedings of the 22nd Conference on Computational Natural Language Learning</i> , pages 65–74, Brussels, Belgium. Association for Computational Linguistics.	733 734 735 736 737 738
	Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. <a href="#">Stress test evaluation for natural language inference</a> . In <i>Proceedings of the 27th International Conference on Computational Linguistics</i> , pages 2340–2353, Santa Fe, New Mexico, USA. Association for Computational Linguistics.	739 740 741 742 743 744 745
	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> .	746 747 748 749
	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.	750 751 752 753

754	Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal,	for robust natural language inference. In <i>Proceed-</i>	810
755	Jason Weston, and Douwe Kiela. 2020. Adversarial	<i>ings of the AAAI Conference on Artificial Intelligence</i> ,	811
756	NLI: A new benchmark for natural language under-	volume 36, pages 11349–11357.	812
757	standing. In <i>ACL</i> , pages 4885–4901. Association for		
758	Computational Linguistics.		
759	Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	813
760	Millican, Jordan Hoffmann, Francis Song, John	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	814
761	Aslanides, Sarah Henderson, Roman Ring, Susan-	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	815
762	nah Young, et al. 2021. Scaling language models:	Bhosale, et al. 2023. Llama 2: Open founda-	816
763	Methods, analysis & insights from training gopher.	tion and fine-tuned chat models. <i>arXiv preprint</i>	817
764	<i>arXiv preprint arXiv:2112.11446</i> .	<i>arXiv:2307.09288</i> .	818
765	Nazneen Fatema Rajani, Bryan McCann, Caiming	Thinh Hung Truong, Yulia Otmakhova, Timothy Bald-	819
766	Xiong, and Richard Socher. 2019. <a href="#">Explain your-</a>	win, Trevor Cohn, Jey Han Lau, and Karin Verspoor.	820
767	<a href="#">self! leveraging language models for commonsense</a>	2022. <a href="#">Not another negation benchmark: The NaN-</a>	821
768	<a href="#">reasoning</a> . In <i>Proceedings of the 57th Annual Meet-</i>	<a href="#">NLI test suite for sub-clausal negation</a> . In <i>Proceed-</i>	822
769	<i>ing of the Association for Computational Linguistics</i> ,	<i>ings of the 2nd Conference of the Asia-Pacific Chap-</i>	823
770	pages 4932–4942, Florence, Italy. Association for	<i>ter of the Association for Computational Linguistics</i>	824
771	Computational Linguistics.	<i>and the 12th International Joint Conference on Natu-</i>	825
772	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-</a>	<i>ral Language Processing (Volume 1: Long Papers)</i> ,	826
773	<a href="#">BERT: Sentence embeddings using Siamese BERT-</a>	pages 883–894, Online only. Association for Compu-	827
774	<a href="#">networks</a> . In <i>Proceedings of the 2019 Conference on</i>	tational Linguistics.	828
775	<i>Empirical Methods in Natural Language Processing</i>	Lewis Tunstall, Edward Beeching, Nathan Lambert,	829
776	<i>and the 9th International Joint Conference on Natu-</i>	Nazneen Rajani, Kashif Rasul, Younes Belkada,	830
777	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages	Shengyi Huang, Leandro von Werra, Clémentine	831
778	3982–3992, Hong Kong, China. Association for Com-	Fourrier, Nathan Habib, et al. 2023. Zephyr: Di-	832
779	putational Linguistics.	rect distillation of lm alignment. <i>arXiv preprint</i>	833
780	Marco Tulio Ribeiro, Sameer Singh, and Carlos	<i>arXiv:2310.16944</i> .	834
781	Guestrin. 2016. <a href="#">"why should i trust you?": Explain-</a>	Alex Wang, Amanpreet Singh, Julian Michael, Felix	835
782	<a href="#">ing the predictions of any classifier</a> . In <i>Proceedings</i>	Hill, Omer Levy, and Samuel Bowman. 2018. <a href="#">GLUE:</a>	836
783	<i>of the 22nd ACM SIGKDD International Conference</i>	<a href="#">A multi-task benchmark and analysis platform for</a>	837
784	<i>on Knowledge Discovery and Data Mining, KDD '16</i> ,	<a href="#">natural language understanding</a> . In <i>Proceedings of the</i>	838
785	page 1135–1144, New York, NY, USA. Association	<i>2018 EMNLP Workshop BlackboxNLP: Analyzing</i>	839
786	for Computing Machinery.	<i>and Interpreting Neural Networks for NLP</i> , pages	840
787	Oscar Sainz, Jon Ander Campos, Iker García-Ferrero,	353–355, Brussels, Belgium. Association for Com-	841
788	Julen Etxaniz, and Eneko Agirre. 2023. <a href="#">Did chatgpt</a>	putational Linguistics.	842
789	<a href="#">cheat on your test?</a>	Jiongxiao Wang, Zichen Liu, Keun Hee Park, Muhao	843
790	William A. Scott. 1962. <a href="#">Cognitive complexity and cog-</a>	Chen, and Chaowei Xiao. 2023a. Adversarial demon-	844
791	<a href="#">nitive flexibility</a> . <i>Sociometry</i> , 25(4):405–414.	stration attacks on large language models. <i>arXiv</i>	845
792	Sofia Serrano and Noah A. Smith. 2019. <a href="#">Is attention in-</a>	<i>preprint arXiv:2305.14950</i> .	846
793	<a href="#">terpretable?</a> In <i>Proceedings of the 57th Annual Meet-</i>	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le,	847
794	<i>ing of the Association for Computational Linguistics</i> ,	Ed H. Chi, Sharan Narang, Aakanksha Chowdhery,	848
795	pages 2931–2951, Florence, Italy. Association for	and Denny Zhou. 2023b. <a href="#">Self-consistency improves</a>	849
796	Computational Linguistics.	<a href="#">chain of thought reasoning in language models</a> . In	850
797	K. Sparck Jones, S. Walker, and S.E. Robertson. 2000.	<i>The Eleventh International Conference on Learning</i>	851
798	<a href="#">A probabilistic model of information retrieval: devel-</a>	<i>Representations</i> .	852
799	<a href="#">opment and comparative experiments: Part 1</a> . <i>Inform-</i>	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,	853
800	<i>ation Processing and Management</i> , 36(6):779–808.	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,	854
801	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao,	Maarten Bosma, Denny Zhou, Donald Metzler, Ed H.	855
802	Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch,	Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy	856
803	Adam R Brown, Adam Santoro, Aditya Gupta,	Liang, Jeff Dean, and William Fedus. 2022a. <a href="#">Emer-</a>	857
804	Adrià Garriga-Alonso, et al. 2022. Beyond the	<a href="#">gent abilities of large language models</a> . <i>Transactions</i>	858
805	imitation game: Quantifying and extrapolating the	<i>on Machine Learning Research</i> . Survey Certifica-	859
806	capabilities of language models. <i>arXiv preprint</i>	tion.	860
807	<i>arXiv:2206.04615</i> .	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	861
808	Joe Stacey, Yonatan Belinkov, and Marek Rei. 2022. Su-	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,	862
809	pervising model attention with human explanations	and Denny Zhou. 2022b. Chain-of-thought prompt-	863
		ing elicits reasoning in large language models. In	864
		<i>NeurIPS</i> .	865



866	Sarah Wiegrefe and Ana Marasovic. 2021. Teach me to explain: A review of datasets for explainable natural language processing. <i>35th Conference on Neural Information Processing Systems (NeurIPS) Track on Datasets and Benchmarks</i> .	<i>of Machine Learning Research</i> , pages 12697–12706. PMLR.	923
867			924
868			
869		Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. 2023. <a href="#">Least-to-most prompting enables complex reasoning in large language models</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	925
870			926
871	Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. <a href="#">A broad-coverage challenge corpus for sentence understanding through inference</a> . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.		927
872			928
873			929
874			930
875			931
876			
877			
878			
879			
880	Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2021. <a href="#">Polyjuice: Generating counterfactuals for explaining, evaluating, and improving models</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6707–6723, Online. Association for Computational Linguistics.		
881			
882			
883			
884			
885			
886			
887			
888			
889	Yuxiang Wu, Matt Gardner, Pontus Stenetorp, and Pradeep Dasigi. 2022. <a href="#">Generating data to mitigate spurious correlations in natural language inference datasets</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2660–2676, Dublin, Ireland. Association for Computational Linguistics.		
890			
891			
892			
893			
894			
895			
896	Jiacheng Ye, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. 2023. <a href="#">Compositional exemplars for in-context learning</a> . In <i>Proceedings of the 40th International Conference on Machine Learning, ICML’23</i> . JMLR.org.		
897			
898			
899			
900			
901	Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. <a href="#">From recognition to cognition: Visual commonsense reasoning</a> . In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> .		
902			
903			
904			
905			
906	Tianyi Zhang, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. <a href="#">Bertscore: Evaluating text generation with bert</a> . In <i>International Conference on Learning Representations</i> .		
907			
908			
909			
910	Yuan Zhang, Jason Baldridge, and Luheng He. 2019. <a href="#">PAWS: Paraphrase adversaries from word scrambling</a> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 1298–1308, Minneapolis, Minnesota. Association for Computational Linguistics.		
911			
912			
913			
914			
915			
916			
917			
918	Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. <a href="#">Calibrate before use: Improving few-shot performance of language models</a> . In <i>Proceedings of the 38th International Conference on Machine Learning</i> , volume 139 of <i>Proceedings</i>		
919			
920			
921			
922			



932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979

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

980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026

<b>Meta-prompt for zero-shot generation</b>
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:
<b>Meta-prompt for few-shot generation</b>
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 4: Meta-prompts used to generate NLEs via ChatGPT in zero- and few-shot scenarios for natural language inference tasks.

<b>Meta-prompt for zero-shot generation</b>
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:
<b>Meta-prompt for few-shot generation</b>
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 5: 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<sub>few</sub>

Tasks	NLEs		
	Vicuna <sub>few</sub>	LLaMA2 <sub>few</sub>	ChatGPT <sub>few</sub>
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)	<b>63.5</b> ( +6.9)

Table 6: ICL performance of Vicuna using (1) standard ICL without NLEs, (2) X-ICL with Vicuna-generated NLEs in a few-shot scenario: Vicuna<sub>few</sub>, (3) X-ICL with LLaMA2-generated NLEs in a few-shot scenario: LLaMA2<sub>few</sub>, (4) X-ICL with ChatGPT-generated NLEs in a few-shot scenario: ChatGPT<sub>few</sub>. Numbers in the parentheses represent differences compared to X-ICL (Human).

and LLaMA2<sub>few</sub>, respectively. To ensure a fair comparison, we employ Vicuna as the underlying model to evaluate X-ICL(Vicuna<sub>few</sub>), X-ICL (LLaMA2<sub>few</sub>), and X-ICL (ChatGPT<sub>few</sub>) on all studied datasets.

Our results, detailed in Table 6, highlight that X-ICL generally gains greater benefit from LLM-generated NLEs as opposed to those produced by humans. Meanwhile, X-ICL (ChatGPT<sub>few</sub>) consistently outperforms X-ICL(Vicuna<sub>few</sub>) and X-ICL (LLaMA2<sub>few</sub>) considerably, except for ANLI R1 and R2. These findings suggest that in order to fully harness the potential of AI-generated NLEs, the employment of a powerful LLM is integral.

**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 X-ICL (ChatGPT<sub>few</sub>), utilizing GPT3.5-turbo as the foundational model. A total of 280 randomly sampled, correctly predicted examples from X-ICL (ChatGPT<sub>few</sub>) 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.
- **insufficient justification**: The NLE requires more support for the prediction.
- **too verbose**: The NLE is overly detailed and includes unnecessary information.

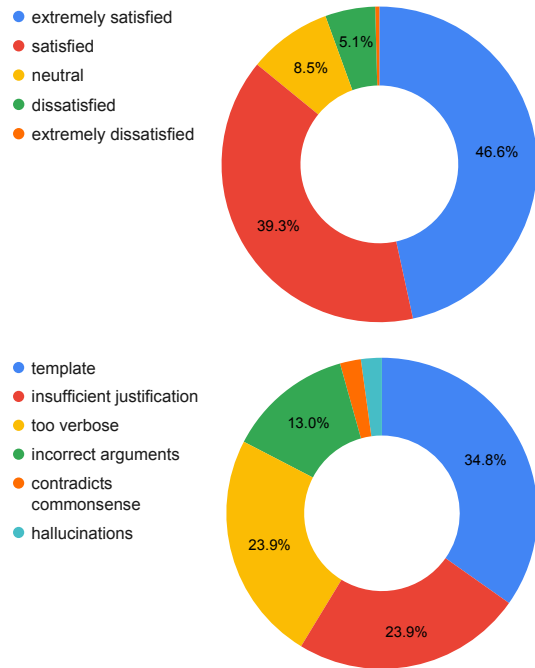


Figure 5: Human evaluation on ChatGPT-generated NLEs for the correct predictions from X-ICL (ChatGPT<sub>few</sub>). **Top**: distribution of satisfaction scores. **Bottom**: distribution of reasons for dissatisfaction.

- **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 5, 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.