# PromptRobust: Towards Evaluating the Robustness of Large Language Models on Adversarial Prompts

**Kaijie Zhu**[1,2][*] **Jindong Wang**[1][†] **Jiaheng Zhou**[2] **, Zeek Wang**[1] **, Hao Chen**[3] **, Yidong Wang**[4] **,**
**Linyi Yang**[5] **, Wei Ye**[4] **, Neil Zhenqiang Gong**[6] **, Yue Zhang**[5] **, Xing Xie**[1]
[1]Microsoft Research  [2]Institute of Automation, CAS  [3]Carnegie Mellon University
[4]Peking University  [5]Westlake University  [6]Duke University

## Abstract

The increasing reliance on Large Language Models (LLMs) necessitates a comprehensive understanding of their robustness to prompts. In this paper, we introduce PromptRobust, a robustness benchmark designed to measure LLMs' resilience to adversarial prompts. This study uses a plethora of adversarial textual attacks on prompts across multiple levels: character, word, sentence, and semantic. The adversarial prompts, crafted to mimic plausible user errors like typos or synonyms, aim to evaluate how slight deviations can affect LLM outcomes while maintaining semantic integrity. These prompts are then used in various tasks, including sentiment analysis, natural language inference, reading comprehension, machine translation, and math. We generate $4,788$ adversarial prompts and evaluated over 8 tasks and 13 datasets. Our findings demonstrate that LLMs are not robust to adversarial prompts. Furthermore, we present a comprehensive analysis to understand the mystery behind prompt robustness and its transferability. We then offer insightful analysis and pragmatic recommendations for prompt composition, beneficial to both researchers and everyday users.

## 1 Introduction

Large language models (LLMs) have gained increasing popularity due to their unprecedented performance in various tasks such as sentiment analysis (Wang et al., 2019), question answering (Wang et al., 2019), logical reasoning (Liu et al., 2023a), etc. An *input* to an LLM is the concatenation of a *prompt* and (optionally) a *sample*, where the prompt aims to instruct the LLM what task to perform and the sample is the data for the task.

Given the popular adoption of LLMs, particularly in the safety-critical and decision-making domains, it becomes essential to examine the *robustness* of LLMs to perturbations in an input. Indeed, existing work (Wang et al., 2021; Nie et al., 2020; Wang et al., 2023b; Zhuo et al., 2023; Yang et al., 2023) has at-
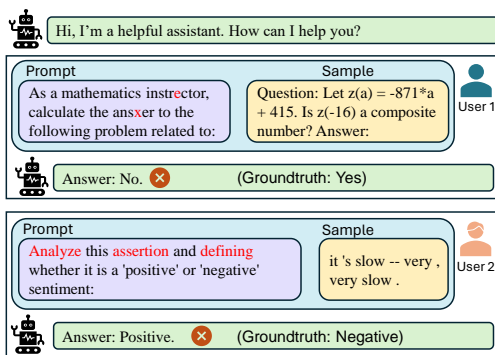


Figure 1: LLMs are not robust to prompts: typos and synonyms lead to errors in math and sentiment analysis problems. The red characters and words are perturbations.

tempted to assess the robustness of LLMs from different perspectives. For instance, AdvGLUE (Wang et al., 2021) and ANLI (Nie et al., 2020) are two public datasets to evaluate the robustness of language models to *adversarial samples*, which are carefully perturbed samples to make a language model produce incorrect responses. In the era of large language models, Wang et al. (2023b) evaluated ChatGPT and other LLMs with respect to their robustness to

---

adversarial samples and out-of-distribution (OOD) samples. Zhuo et al. (2023) evaluated the robustness of LLMs for a particular task called semantic parsing.

These studies demonstrated that the current LLMs are not robust to adversarial and OOD samples for some popular natural language processing tasks. A single prompt is often used to instruct an LLM to perform a task for multiple samples. For example, in a math problem (shown in Figure 1), a prompt can be used for multiple samples (i.e., math problems). Therefore, a perturbed prompt may cause an LLM to output incorrect responses for multiple clean samples. As a result, a perturbed prompt arguably has a larger impact on LLMs than an adversarial sample, as the latter only influences the response of an LLM for a single sample. However, despite its central importance, the robustness of LLMs to prompt perturbations is largely unexplored.

In this paper, we aim to bridge the gap by introducing **PromptRobust**, a comprehensive benchmark designed to evaluate the robustness of LLMs to perturbations in prompts, understanding the factors that contribute to their robustness (or lack thereof), and identifying the key attributes of robust prompts. We consider a variety of prompt perturbations including 1) minor typos, synonyms, and different ways of expressing sentences with the same semantic meaning, which may commonly occur to normal users or developers in their daily use of LLMs in *non-adversarial settings*, as well as 2) perturbations strategically crafted by attackers in *adversarial settings*. With a slight abuse of terminology, we call such a perturbed prompt in both scenarios *adversarial prompt*. Figure 1 shows examples of adversarial prompts with typos and synonyms that lead to incorrect responses.

PromptRobust consists of *prompts*, *attacks*, *models*, *tasks*, *datasets*, and *analysis*. Specifically, we evaluate 4 types of prompts: zero-shot (ZS), few-shot (FS), role-oriented, and task-oriented prompts. We create 4 types of attacks (called *prompt attacks*) to craft adversarial prompts: *character-level*, *word-level*, *sentence-level*, and *semantic-level* attacks by extending 7 adversarial attacks (Li et al., 2019; Gao et al., 2018; Li et al., 2020; Jin et al., 2019; Naik et al., 2018; Ribeiro et al., 2020) that were originally designed to generate adversarial samples. We note that, although we call them attacks, their generated adversarial prompts also serve as testbeds for *mimicking* potential diverse prompts with naturally occurred perturbations from real LLM users. PromptRobust spans across 9 prevalent LLMs, ranging from smaller models such as Flan-T5-large (Chung et al., 2022) to larger ones like ChatGPT (OpenAI, 2023a) and GPT-4 (OpenAI, 2023b). Moreover, we select 8 tasks for evaluation, namely, sentiment analysis (SST-2 (Socher et al., 2013)), grammar correctness (CoLA (Warstadt et al., 2018)), duplicate sentence detection (QQP (Wang et al., 2017) and MRPC (Dolan & Brockett, 2005)), natural language inference (MNLI (Williams et al., 2018), QNLI (Wang et al., 2019), RTE (Wang et al., 2019), and WNLI (Levesque et al., 2012)), multi-task knowledge (MMLU (Hendrycks et al., 2021)), reading comprehension (SQuAD V2 (Rajpurkar et al., 2018)), translation (UN Multi (Eisele & Chen, 2010) and IWSLT 2017 (Cettolo et al., 2017)), and math problem-solving (Mathematics (Saxton et al., 2019)). In total, we created $4,788$ adversarial prompts, representing diverse, practical, and challenging scenarios.

We carry out extensive experiments and analysis using PromptRobust. The results highlight a prevailing lack of robustness to adversarial prompts among current LLMs, with word-level attacks proving the most effective (39% average performance drop in all tasks). We delve into the reasons behind this vulnerability by exploring LLMs' attention weights of each word in inputs for erroneous responses associated with clean and adversarial inputs, where an adversarial input is the concatenation of an adversarial prompt and a clean sample. Our findings reveal that adversarial prompts cause LLMs to shift their focus towards the perturbed elements thus producing wrong responses. We also examine the transferability of adversarial prompts between models, and suggest a successful transferability of adversarial prompts from one LLM to another. Furthermore, we analyze word frequency patterns to guide future research in improving robustness and to aid end-users in crafting more robust prompts. We conclude by discussing potential strategies for improving robustness.

To summarize, our contributions are as follows:

1. We introduce PromptRobust, the *first* systematic benchmark for evaluating, understanding, and analyzing the robustness of LLMs to adversarial prompts.

2. We conduct comprehensive evaluations on the robustness of LLMs to adversarial prompts and perform extensive analysis, including visual explanations for observed vulnerabilities, transferability analysis of adversarial prompts, and word frequency analysis to offer practical guidance to downstream users and prompt engineers to craft more robust prompts.

## 2 PromptRobust

### 2.1 Prompts and models

We investigate four different types of prompt. **Task-oriented** prompts explicitly describe the task the model is required to perform, which encourages the model to generate task-specific outputs based solely on its pre-training knowledge. While **role-oriented** prompts typically frame the model as an entity with a specific role, such as an expert, advisor, or translator. By incorporating role information, these prompts aim to implicitly convey the expected output format and behavior. Each of the two categories of prompts can be designed for both **zero-shot (ZS)** and **few-shot (FS)** learning scenarios. In the zero-shot scenario, an input is defined as $[P, x]$, where $P$ denotes a prompt, $x$ is a sample, and $[,]$ denotes the concatenation operation. For the few-shot scenario, some examples are added to the input, resulting in the format $[P, E, x]$, where $E$ represents the examples. For instance, $E = \{[x_1, y_1], [x_2, y_2], [x_3, y_3]\}$ represents three examples in a three-shot learning scenario. In our experiments, we randomly select three examples in the training set of a task and append them to a prompt. Appendix A.1 shows examples of different types of prompts.

Our evaluation includes a diverse set of LLMs to comprehensively assess their performance across various tasks and domains. The models we consider are as follows: Flan-T5-large (Chung et al., 2022) (0.8B), Dolly-6B (Databricks, 2023), Vicuna-13B (Chiang et al., 2023), Llama2-13b-chat (Touvron et al., 2023b), Cerebras-GPT-13B (Dey et al., 2023), GPT-NEOX-20B (Black et al., 2022), Flan-UL2 (20B) (Brain, 2023), ChatGPT (OpenAI, 2023a), and GPT-4 (OpenAI, 2023b).[1] By incorporating LLMs with different architectures and sizes, we aim to provide insights into their strengths and weaknesses, ultimately facilitating model selection for a specific application or use case. Details of these LLMs are in Appendix B.1.

### 2.2 Attacks

Given a single sample $x$ and its label $y$, a textual adversarial attack aims to find a perturbation $\delta$ such that an LLM $f_\theta$ produces an incorrect response. Formally, $\delta$ is found by solving the following optimization problem: $\max_{\delta \in \mathcal{C}} \mathcal{L}[f_\theta(x + \delta); y]$, where $x + \delta$ is the adversarial sample, $f_\theta(x + \delta)$ is the response of the LLM when taking the adversarial sample alone as input, $\mathcal{C}$ indicates the constraints for the perturbation $\delta$, and $\mathcal{L}$ represents a loss function.

**Prompt attack.** In this paper, our focus is to attack the *prompts* rather than samples. This is due to the popularity of LLMs in different applications, which generate responses using in-context learning on prompts (i.e., instructions) and samples. Prompts are either input by users or generated by the system or developers. Moreover, the ultimate purpose of performing such "attack" is actually *not* to genuinely attack the models, but to *simulate* possible perturbations that may naturally occur in real situations. Table 8 shows multiple prompts generated by adversarial approaches that are used to mimic possible user prompts, which are popular errors or expressions made by users. Since users can make different mistakes when entering prompts, such as typos, different word usage, different sentence styles, etc., the study on prompt robustness is necessary to understand LLMs.

We denote an input to LLMs as $[P, x]$, where $P$ is a prompt, $x$ is a sample, and $[,]$ denotes concatenation. Note that in the few-shot learning scenario, a few examples are appended to the prompt; and the sample $x$ is optional in certain application scenarios. Our prompt attack can also be extended to such scenarios, but we use the notation $[P, x]$ for simplicity. Given a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i \in [N]}$ with $N$ samples and their ground-truth labels, a prompt

---

[1]We did not perform prompt attacks on GPT-4 by optimizing the adversarial algorithms since it requires massive rounds of communications and is too costly. We used the adversarial prompts generated by ChatGPT to evaluate GPT-4 since the adversarial prompts can be transferred (Sec. 4.4).

attack aims to perturb $P$ such that an LLM $f_\theta$ produces incorrect responses for all samples in the dataset $\mathcal{D}$.

**Definition 2.1** (Prompt Attack)**.** Given an LLM $f_\theta$, a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i \in [N]}$, and a clean prompt $P$, the objective of a prompt attack can be formulated as follows:

$$\max_{\delta \in \mathcal{C}} \sum_{(x;y) \in \mathcal{D}} \mathcal{L}[f_\theta([P + \delta, x]), y], \tag{1}$$

where $\delta$ is the textual perturbation added to the clean prompt $P$ and $\mathcal{C}$ is the allowable perturbation set, i.e., perturbation constraint. We note that this attack is analogous to universal adversarial perturbation (UAP) (Moosavi-Dezfooli et al., 2017; Brown et al., 2017) and universal adversarial trigger (UAT) (Wallace et al., 2019), extending these concepts to the realm of prompts. Appendix A.4 shows more comparisons.

**Different attacks.** We then modify the existing black-box textual attacks to implement Eq. (1) due to their efficiency and no reliance on the model gradient. Thus, both open-sourced and proprietary LLMs can be attack targets. Our instantiations span four distinct levels, capturing a broad spectrum of complexities from simple character manipulations to sophisticated semantic alterations. Details of each attack are in Appendix A.3.

- **Character-level:** We employ TextBugger (Li et al., 2019) and DeepWordBug (Gao et al., 2018), which manipulate texts by introducing typos or errors to words, e.g., by adding, deleting, repeating, replacing, and permuting characters for certain words.
- **Word-level:** We use BertAttack (Li et al., 2020) and TextFooler (Jin et al., 2019) to replace words with synonyms or contextually similar words to deceive LLMs.
- **Sentence-level:** We implement StressTest (Naik et al., 2018) and CheckList (Ribeiro et al., 2020) to append irrelevant or extraneous sentences to the end of prompts, intending to distract LLMs. For StressTest, we adopt settings similar to those in (Wang et al., 2019), appending "and true is true", "and false is not true", or "and true is true" for five times to the end of a prompt. For the CheckList attack, we generate 50 random sequences consisting of alphabets and digits, each with a length of 10, and append this random sequence to the end of a prompt.
- **Semantic-level:** We simulate the linguistic behavior of people from different countries by choosing 6 common languages (Chinese, French, Arabic, Spanish, Japanese, and Korean) and constructing 10 prompts for each language per dataset. These prompts are then translated into English, introducing linguistic nuances and variations that could potentially impact LLMs.

**Semantic preservation of adversarial prompts.**

Are adversarial prompts realistic? Our purpose is to simulate plausible user errors; thus, it is imperative that these prompts preserve semantic integrity, ensuring they remain both *acceptable* and *imperceptible* to human comprehension. It is of paramount importance that our adversarially engineered prompts retain coherence and realism. Our human study in Appendix A.2 shows that at least 85% of subjects agreed that the generated prompts are acceptable.

Table 1: Statistics of datasets used in this paper.

| Task | Dataset | #Sample | #Class | #[Adv. prompt, sample] |
|------|---------|---------|--------|------------------------|
| Sentiment analysis | SST2 | 872 | 2 | 73,248 |
| Grammar correctness | CoLA | 1,000 | 2 | 84,000 |
| Duplicate sentence detection | QQP | 1,000 | 2 | 84,000 |
| | MRPC | 408 | 2 | 34,272 |
| Natural language inference | MNLI | 1,000 | 3 | 84,000 |
| | QNLI | 1,000 | 2 | 84,000 |
| | RTE | 277 | 2 | 23,268 |
| | WNLI | 71 | 2 | 5,964 |
| Multi-task knowledge | MMLU | 564 | 4 | 47,376 |
| Reading comprehension | SQuAD V2 | 200 | - | 16,800 |
| Translation | Multi UN | 99 | - | 8,316 |
| | IWSLT 2017 | 100 | - | 8,400 |
| Math reasoning | Math | 160 | - | 13,440 |

## 2.3 Tasks and datasets

Currently, PromptRobust supports 8 tasks and 13 datasets across sentiment analysis to math reasoning. Due to space limit, we leave the details of the datasets in Appendix B.2.

## 3 Experiments

**Setup.** The extensive computational need to generate an adversarial prompt requires iterating throughout the dataset 100 on average. Therefore, the evaluation of an entire

Table 2: The APDR and standard deviations of different attacks on different datasets.

| Dataset | Character-level | | Word-level | | Sentence-level | | Semantic-level |
|---|---|---|---|---|---|---|---|
| | TextBugger | DeepWordBug | TextFooler | BertAttack | CheckList | StressTest | Semantic |
| SST-2 | 0.25±0.39 | 0.18±0.33 | 0.35±0.41 | 0.34±0.44 | 0.22±0.36 | 0.15±0.31 | 0.28±0.35 |
| CoLA | 0.39±0.40 | 0.27±0.32 | 0.43±0.35 | 0.45±0.38 | 0.23±0.30 | 0.18±0.25 | 0.34±0.37 |
| QQP | 0.30±0.38 | 0.22±0.31 | 0.31±0.36 | 0.33±0.33 | 0.18±0.30 | 0.06±0.26 | 0.40±0.39 |
| MRPC | 0.37±0.42 | 0.34±0.41 | 0.37±0.41 | 0.42±0.38 | 0.24±0.37 | 0.25±0.33 | 0.39±0.39 |
| MNLI | 0.32±0.40 | 0.18±0.29 | 0.32±0.39 | 0.34±0.36 | 0.14±0.24 | 0.10±0.25 | 0.22±0.24 |
| QNLI | 0.38±0.39 | 0.40±0.35 | 0.50±0.39 | 0.52±0.39 | 0.25±0.39 | 0.23±0.33 | 0.40±0.35 |
| RTE | 0.33±0.41 | 0.25±0.35 | 0.37±0.44 | 0.40±0.42 | 0.18±0.32 | 0.17±0.24 | 0.42±0.40 |
| WNLI | 0.39±0.42 | 0.31±0.37 | 0.41±0.43 | 0.41±0.40 | 0.24±0.32 | 0.20±0.27 | 0.49±0.39 |
| MMLU | 0.21±0.24 | 0.12±0.16 | 0.21±0.20 | 0.40±0.30 | 0.13±0.18 | 0.03±0.15 | 0.20±0.19 |
| SQuAD V2 | 0.09±0.17 | 0.05±0.08 | 0.25±0.29 | 0.31±0.32 | 0.02±0.03 | 0.02±0.04 | 0.08±0.09 |
| IWSLT | 0.08±0.14 | 0.10±0.12 | 0.27±0.30 | 0.12±0.18 | 0.10±0.10 | 0.17±0.19 | 0.18±0.14 |
| UN Multi | 0.06±0.08 | 0.08±0.12 | 0.15±0.19 | 0.10±0.16 | 0.06±0.07 | 0.09±0.11 | 0.15±0.18 |
| Math | 0.18±0.17 | 0.14±0.13 | 0.49±0.36 | 0.42±0.32 | 0.15±0.11 | 0.13±0.08 | 0.23±0.13 |
| Avg | 0.21±0.30 | 0.17±0.26 | 0.31±0.33 | 0.33±0.34 | 0.12±0.23 | 0.11±0.23 | 0.22±0.26 |

Table 3: The APDR on different LLMs.

| Dataset | T5-large | Vicuna | Llama2 | UL2 | ChatGPT | GPT-4 |
|---|---|---|---|---|---|---|
| SST-2 | 0.04±0.11 | 0.83±0.26 | 0.24±0.33 | 0.03±0.12 | 0.17±0.29 | 0.24±0.38 |
| CoLA | 0.16±0.19 | 0.81±0.22 | 0.38±0.32 | 0.13±0.20 | 0.21±0.31 | 0.13±0.23 |
| QQP | 0.09±0.15 | 0.51±0.41 | 0.59±0.33 | 0.02±0.04 | 0.16±0.30 | 0.16±0.38 |
| MRPC | 0.17±0.26 | 0.52±0.40 | 0.84±0.27 | 0.06±0.10 | 0.22±0.29 | 0.04±0.06 |
| MNLI | 0.08±0.13 | 0.67±0.38 | 0.32±0.32 | 0.06±0.12 | 0.13±0.18 | -0.03±0.02 |
| QNLI | 0.33±0.25 | 0.87±0.19 | 0.51±0.39 | 0.05±0.11 | 0.25±0.31 | 0.05±0.23 |
| RTE | 0.08±0.13 | 0.78±0.23 | 0.68±0.39 | 0.02±0.04 | 0.09±0.13 | 0.03±0.05 |
| WNLI | 0.13±0.14 | 0.78±0.27 | 0.73±0.37 | 0.04±0.03 | 0.14±0.12 | 0.04±0.04 |
| MMLU | 0.11±0.18 | 0.41±0.24 | 0.28±0.24 | 0.05±0.11 | 0.14±0.18 | 0.04±0.04 |
| SQuAD V2 | 0.05±0.12 | - | - | 0.10±0.18 | 0.22±0.28 | 0.27±0.31 |
| IWSLT | 0.14±0.17 | - | - | 0.15±0.11 | 0.17±0.26 | 0.07±0.14 |
| UN Multi | 0.13±0.14 | - | - | 0.05±0.05 | 0.12±0.18 | -0.02±0.01 |
| Math | 0.24±0.21 | - | - | 0.21±0.21 | 0.33±0.31 | 0.02±0.18 |
| Avg | 0.13±0.19 | 0.69±0.34 | 0.51±0.39 | 0.08±0.14 | 0.18±0.26 | 0.08±0.21 |

Table 4: APDR on different prompts.

| Dataset | ZS-task | ZS-role | FS-task | FS-role |
|---|---|---|---|---|
| SST-2 | 0.31±0.39 | 0.28±0.35 | 0.22±0.38 | 0.24±0.39 |
| CoLA | 0.43±0.35 | 0.43±0.38 | 0.24±0.28 | 0.25±0.36 |
| QQP | 0.43±0.42 | 0.34±0.43 | 0.16±0.21 | 0.14±0.20 |
| MRPC | 0.44±0.44 | 0.51±0.43 | 0.24±0.32 | 0.23±0.30 |
| MNLI | 0.29±0.35 | 0.26±0.33 | 0.19±0.29 | 0.21±0.33 |
| QNLI | 0.46±0.39 | 0.51±0.40 | 0.30±0.34 | 0.32±0.36 |
| RTE | 0.33±0.39 | 0.35±0.40 | 0.31±0.39 | 0.27±0.38 |
| WNLI | 0.36±0.36 | 0.39±0.39 | 0.37±0.41 | 0.33±0.38 |
| MMLU | 0.25±0.23 | 0.22±0.26 | 0.18±0.23 | 0.14±0.20 |
| SQuAD V2 | 0.16±0.26 | 0.20±0.28 | 0.06±0.11 | 0.07±0.12 |
| IWSLT | 0.18±0.22 | 0.24±0.25 | 0.08±0.09 | 0.11±0.10 |
| UN Multi | 0.17±0.18 | 0.15±0.16 | 0.04±0.07 | 0.04±0.07 |
| Math | 0.33±0.26 | 0.39±0.30 | 0.16±0.18 | 0.17±0.17 |
| Avg | 0.33±0.36 | 0.34±0.37 | 0.21±0.31 | 0.21±0.31 |

dataset using LLMs is unfeasible. To alleviate the computational constraint and preserve a fair study process, we adopt a sampling strategy that entails selecting a subset of samples from the validation or test sets across various datasets. The statistics of each dataset and tasks are summarized in Table 1.[2] The sampling details are in Appendix C.

We initially assess the performance of all LLMs without prompt attacks to provide a performance baseline. We find that certain LLMs do not even demonstrate satisfactory performance with clean prompts, narrowing our selection to 6 LLMs: Flan-T5-large, Vicuna-13B, Llama2-13B-chat, UL2, ChatGPT, and GPT-4. Further details and discussions on clean prompt performance in all LLMs are available in Appendix C.2. We generate 10 distinct prompts for both role-oriented and task-oriented categories. Each prompt can be augmented with three examples, forming the few-shot prompts. In total, we have 40 prompts for each dataset on each LLM. For better efficiency and performance, we select the top 3 best-performing prompts of each type to perform prompt attacks. As a result, we evaluate the adversarial vulnerabilities of 9 LLMs across 13 datasets, encompassing a total of $4,788$ prompts[3] and their respective adversarial counterparts. This comprehensive evaluation allows us to gain valuable insight into the robustness and performance of LLMs in a wide range of scenarios and prompt styles.

**Evaluation metrics.** Considering the diverse evaluation metrics across tasks and varying baseline performances across models and datasets, the absolute performance drop may not provide a meaningful comparison. Thus, we introduce a unified metric, the *Performance Drop Rate* (PDR). PDR quantifies the relative performance decline following a prompt attack, offering a contextually normalized measure to compare different attacks, datasets, and models. The PDR is given by: $PDR(A, P, f_\theta, \mathcal{D}) = 1 - \frac{\sum_{(x;y)\in\mathcal{D}} \mathcal{M}[f_\theta([A(P),x]),y]}{\sum_{(x;y)\in\mathcal{D}} \mathcal{M}[f_\theta([P,x]),y]}$, where $A$ is the

---

[2]In Table 1, the last column denotes the total evaluation sample size for each dataset on each model. For instance, there are 872 test samples in SST-2 dataset and each sample should go through 7 adversarial attacks on 4 types of prompts, each with 3 prompts, thus the test size for each model is $872 \times 7 \times 4 \times 3 = 73248$.

[3]$4,788 = 3 \times 4 \times 5 \times 13 \times 7 - 336 \times 3 \times 2$, where the numbers on the R.H.S. denote #attacked prompts, #prompt types, #LLMs, #datasets, and #attacks, respectively. We did not conduct attacks on Vicuna, Llama2-13B-chat on certain datasets because the output of these datasets are meaningless, so that we subtract $336 \times 3 \times 2$ prompts.

adversarial attack applied to prompt $P$, and $\mathcal{M}[\cdot]$ is the evaluation function: for classification task, $\mathcal{M}[\cdot]$ is the indicator function $\mathbb{1}[\hat{y}, y]$ which equals to 1 when $\hat{y} = y$, and 0 otherwise; for reading comprehension task, $\mathcal{M}[\cdot]$ is the F1-score; for translation tasks, $\mathcal{M}[\cdot]$ is the Bleu metric (Papineni et al., 2002). Note that a negative PDR implies that adversarial prompts can enhance performance.

## 4 Results and analysis

### 4.1 Results across attacks, models, and prompts

We report and discuss the Average PDR (APDR) across different attacks, LLMs, and prompts. Our main results are based on *all* the prompts, whose conclusions are in consistent with Appendix C.4 where we show result by excluding unacceptable adversarial prompts. Note that although our semantic preserving study demonstrated that at least 85% of adversarial prompts are acceptable, there are still some adversarial prompts diverged from their intended semantic meaning. Furthermore, note that the discrepancies in APDR variance values are due to varying PDR values across different attacks, prompts, models and datasets.

**Analysis on attacks.** Table 2 summarizes the APDR of 7 attacks on 13 datasets, calculated by $APDR_A(A, \mathcal{D}) = \frac{1}{|\mathcal{P}|} \frac{1}{|\mathcal{F}|} \sum_{P \in \mathcal{P}} \sum_{f_\theta \in \mathcal{F}} PDR(A, P, f_\theta, \mathcal{D})$, where $\mathcal{P}$ is the set of 4 types of prompts and $\mathcal{F}$ is the set of all models. The results offer several key insights. Firstly, attack effectiveness is highly variable, with word-level attacks proving the most potent, leading to an average performance decline of 33% across *all* datasets. Character-level attacks rank second, inducing a 20% performance drop in most datasets. Notably, semantic-level attacks exhibit potency nearly commensurate with character-level attacks, emphasizing the profound impact of nuanced linguistic variations on LLMs' performance. On the contrary, sentence-level attacks pose less of a threat, suggesting adversarial interventions at this level have a diminished effect. Moreover, the effect of prompt attack varies across different datasets. For instance, StressTest attacks on SQUAD V2 yield a mere 2% performance drop, while inflicting a 25% drop on MRPC. Furthermore, we observe that the StressTest attack paradoxically bolsters the model's performance in some datasets, we delve into this phenomenon in Appendix D.3.

Note that while character-level attacks are detectable by grammar detection tools, word- and semantic-level attacks underscore the importance of robust semantic understanding and accurate task presentation/translation for LLMs. A comprehensive understanding of these nuances will inform a deeper comprehension of adversarial attacks on LLMs.

**Analysis on LLMs.** Table 3 summarizes the APDR of 9 LLMs on 13 datasets, calculated by $APDR_{f_\theta}(f_\theta, \mathcal{D}) = \frac{1}{|\mathcal{A}|} \frac{1}{|\mathcal{P}|} \sum_{A \in \mathcal{A}} \sum_{P \in \mathcal{P}} PDR(A, P, f_\theta, \mathcal{D})$, where $\mathcal{P}$ is the set of 4 types of prompts and $\mathcal{A}$ is the set of 7 attacks. Our analysis reveals that GPT-4 and UL2 significantly outperform other models in terms of robustness, followed by T5-large, ChatGPT, and Llama2, with Vicuna presenting the least robustness. The robustness against adversarial prompts of UL2, T5-large, and ChatGPT varies across datasets, with UL2 and T5-large showing less vulnerability to attacks on sentiment classification (SST-2), most NLI tasks, and reading comprehension (SQuAD V2). Specifically, UL2 excels in translation tasks, while ChatGPT displays robustness in certain NLI tasks. Vicuna, however, exhibits consistently high susceptibility to attacks across all tasks. It can be seen that, given the same adversarial prompts generated by ChatGPT, GPT-4 exhibits superior robustness in all tasks. However, it is crucial to realize that this observed robustness might attribute to the weak transferability of the adversarial prompts crafted specifically for ChatGPT. In the future, the performance of GPT-4 and ChatGPT could be improved since these proprietary models keep evolving. **Models vs. attacks.** We show in Table 10 the relation between models and attacks. Generally, word-level attacks emerge as the most potent, and BertAttack consistently outperforms others across all models. However, no discernible pattern emerges for the efficacy of the other attacks. For instance, while TextBugger proves more effective than DeepWordBug for some models such as Llama2 and ChatGPT, the inverse holds true for T5-large. Notably, Vicuna and Llama2 are distinctly vulnerable to sentence-level attacks, in contrast to models like T5-large and ChatGPT, which remain largely unaffected. Such observations may hint at inherent vulnerabilities specific to Llama-based models.

Table 5: Attention visualization of samples that are *correctly classified by clean prompts but misclassified by adv. prompts*. For each attack, the above is the *clean prompt* with sample text, the below is the corresponding *adversarial prompt* with the same sample text. N=Negative, P=Positive and N/A means the response is not available. The green and red color denote right and wrong answers, respectively. Color intensity denotes different attention weights (heavier color means larger weights).

| Attack | Pred. | [Prompt, sample] |
|--------|-------|------------------|
| BertAttack | N | In the role of a sentiment analysis tool, respond with 'positive' or 'negative' to classify this statement:the title not only describes its main characters , but the lazy people behind the camera as well . Answer: |
| | P | how the role of a compliment analysis tool, responses with 'positive' or 'negative' to mood this statement:the title not only describes its main characters , but the lazy people behind the camera as well . Answer: |
| CheckList | P | Given the context of this text, indicate if the emotion conveyed is 'positive' or 'negative':good film , but very glum . Answer: |
| | N | Given the context of this text, indicate if the emotion conveyed is 'positive' or 'negative' LkF0FZxMZ4:good film , but very glum . Answer: |
| DeepWordBug | N | Serving as a sentiment evaluation model, determine if the given statement is 'positive' or negative . Classify:i had to look away - this was god awful . Answer: |
| | N/A | Servign as a sentimBnt envaluation model, Determine if the liven statemen is 'positive' or 'negative'. Classhfy: i had to look away - this was god awful . Answer: |
| Semantic | N | In the role of a sentiment analysis tool, respond with 'positive' or 'negative' to classify this statement:bad . Answer: |
| | P | Classify what you're trying to convey in this sentence as 'positive' if it's positive, and 'negative' if it's negative:bad . Answer: |

Table 6: The APDR of transferability of several LLMs.

| Attacks | Chat → T5 | Chat → UL2 | Chat → V | T5 → Chat | T5 → UL2 | T5 → V | UL2 → Chat | UL2 → T5 | UL2 → V | V → Chat | V → T5 | V → UL2 |
|---------|-----------|------------|----------|-----------|----------|--------|------------|----------|---------|----------|--------|---------|
| BertAttack | 0.05±0.17 | 0.08±0.19 | 0.08±0.88 | 0.18±0.32 | 0.11±0.23 | -1.39±5.67 | 0.15±0.27 | 0.05±0.11 | -0.70±3.18 | 0.06±0.19 | 0.05±0.11 | 0.03±0.12 |
| CheckList | 0.00±0.04 | 0.01±0.03 | 0.19±0.39 | 0.00±0.07 | 0.01±0.03 | -0.09±0.64 | 0.01±0.06 | 0.01±0.04 | -0.13±1.80 | -0.01±0.04 | 0.00±0.01 | 0.00±0.00 |
| TextFooler | 0.04±0.08 | 0.03±0.09 | -0.25±1.03 | 0.11±0.23 | 0.08±0.16 | -0.30±2.09 | 0.11±0.21 | 0.07±0.18 | -0.17±1.46 | 0.04±0.16 | 0.02±0.06 | 0.00±0.01 |
| TextBugger | -0.00±0.09 | -0.01±0.05 | 0.02±0.94 | 0.04±0.15 | 0.01±0.04 | -0.45±3.43 | 0.04±0.13 | 0.02±0.07 | -0.84±4.42 | 0.03±0.13 | 0.01±0.05 | 0.00±0.01 |
| DeepWordBug | 0.03±0.11 | 0.01±0.03 | 0.10±0.46 | 0.00±0.06 | 0.01±0.02 | -0.18±1.20 | 0.01±0.10 | 0.02±0.06 | -0.09±0.75 | 0.00±0.03 | 0.02±0.11 | 0.00±0.01 |
| StressTest | 0.04±0.17 | 0.03±0.10 | 0.01±0.48 | -0.01±0.06 | 0.03±0.06 | 0.04±0.80 | 0.00±0.04 | 0.05±0.16 | 0.06±0.45 | 0.00±0.04 | 0.09±0.18 | 0.02±0.08 |
| Semantic | 0.04±0.12 | 0.02±0.06 | 0.25±0.47 | 0.07±0.27 | 0.00±0.03 | -0.81±4.14 | 0.02±0.11 | -0.13±0.72 | -0.50±1.59 | 0.07±0.11 | 0.00±0.05 | 0.00±0.02 |

**Analysis on types of prompts.** Table 4 summarizes the APDR of 4 types of prompts on 13 datasets, calculated by $APDR_t(\mathcal{D}) = \frac{1}{|\mathcal{A}|} \frac{1}{|\mathcal{P}_t|} \frac{1}{|\mathcal{F}|} \sum_{\mathcal{A} \in \mathcal{A}} \sum_{P \in \mathcal{P}_t} \sum_{f_\theta \in \mathcal{F}} PDR(A, P, f_\theta, \mathcal{D})$, where $\mathcal{P}_t$ is the set of prompts of certain type $t$, $\mathcal{A}$ is the set of 7 attacks and $\mathcal{F}$ is the set of all models. In our analysis, few-shot prompts consistently demonstrate superior robustness compared to zero-shot prompts across all datasets. Furthermore, while task-oriented prompts marginally outperform role-oriented prompts in overall robustness, both show varying strengths across different datasets and tasks. For example, role-oriented prompts present increased robustness within the SST-2 and QQP datasets, while task-oriented prompts are more resilient within the MRPC, QNLI, SQuAD V2, and IWSLT datasets. Insights into different effects of prompt types on model vulnerability can inform better prompt design and tuning strategies, enhancing LLMs robustness against adversarial prompt attacks.

## 4.2 Results on model size, fine-tuning, and adversarial inputs

**Model size and fine-tuning.** We analyze the performance on different models sizes using the Llama2 series (7B, 13B, and 70B). Our results in Appendix C.3 show that larger models are typically more robust than smaller ones, but exceptions can occur when smaller models outperform larger ones, which is an interesting finding that can trigger future research. We also evaluated the impact of fine-tuning using vanilla Llama2 and Llama2-chat models in Appendix C.3, indicating that fine-tuned models are generally better at adversarial prompts.

**Attacking both prompts and samples.** We attacked both prompts and tested on Ad-vGLUE (Wang et al., 2021), which contains adversarial samples. Our results in Table 11 show that attacking both will perform even worse. However, intriguing things happen, since attacking both can sometimes enhance the performance that needs further effort.

## 4.3 Understanding the vulnerability of LLMs to adversarial prompts

We study the magic behind adversarial prompts to analyze why they lead to errors for LLMs from different aspects: attention visualization, erroneous response analysis (Appendix D.1), and sentence-level analysis (Appendix D.3).

We visualize the attentions to investigate the influence of adversarial prompts on LLMs' focus on input words. Specifically, we propose two attention visualization techniques: 1) Attention by Gradient, which assigns an attention score to each word based on the gradient norm, and 2) Attention by Deletion, which assigns an attention score to each word by examining the absolute change in loss when the word is removed. The details of these methods can be found in Appendix D.2. Both techniques produce similar results; hence, we focus on results from the Attention by Gradient method for simplicity. Our key findings, as demonstrated in Table 5, are as follows:

- **Clean prompts: efficient attention allocation.** LLMs predominantly concentrate on key terms within clean prompts, aiding in accurate classifications. For instance, for clean prompts of BertAttack in Table 5, LLMs mainly allocate attention to the term 'lazy', correctly deducing a 'Negative' sentiment.
- **Adversarial prompts: attention divergence.** Adversarial prompts can reroute LLMs' attention from integral text segments, causing misclassifications. In some attacks like CheckList and StressTest, the model simultaneously concentrates on the target text and adversarial content, amplifying its susceptibility to adversarial perturbations. For instance, introducing a random sequence 'LKF0FZxMZ4' during a CheckList attack distracts the model, reducing focus on the critical word 'good' for accurate classification. In other attacks, such as BertAttack and DeepWordBug, the model's attention is entirely diverted from the text requiring classification towards adversarial prompts, leading to a significant shift in focus. For example, in DeepWordBug attack, typos in specific words divert the model's attention from 'awful' to the altered word 'Qetermine'.

### 4.4 Transferability of adversarial prompts

Table 6 displays the effectiveness of various attacks in transferring adversarial prompts between several LLMs. For each dataset and prompt type, we selected the most vulnerable prompts generated by a source model (e.g., ChatGPT). These prompts were then utilized to launch transfer attacks against the target models (e.g., T5-large). The impact of these transfer attacks was quantified by calculating $APDR_{\text{transfer}}(A, f_\theta^{\text{target}}) = \frac{1}{|\mathcal{P}_{\text{source}}|} \frac{1}{|\mathbb{D}|} \sum_{P \in \mathcal{P}_{\text{source}}} \sum_{\mathcal{D} \in \mathbb{D}} PDR(A, P, f_\theta^{\text{target}}, \mathcal{D})$, where $f_\theta^{\text{target}}$ is the target model, $\mathcal{P}_{\text{source}}$ is the prompts selected from source model and $\mathbb{D}$ is the set of all datasets.

In general, we observe that adversarial prompts exhibit some degree of transferability. However, it is marginal compared to Table 2 and 3. Specifically, the APDR in the target model by adversarial prompts from source model is small compared to the original APDR of the source model. Furthermore, the standard deviation tends to be larger than the APDR, indicating that the transferability is inconsistent. Some adversarial prompts can be successfully transferred, causing a performance drop, while others may unexpectedly improve the performance of the target model. A prime example is the BertAttack transfer from UL2 to Vicuna, which resulted in a $-0.70(3.18)$ value, suggesting an unanticipated enhancement in Vicuna's performance when subjected to these adversarial prompts. These phenomena illustrate the complex robustness traits of different models. The transferability to ChatGPT is better compared to T5-large and UL2. This suggests an avenue to generate adversarial prompts to attack black-box models such as ChatGPT by training on small models like T5-large, which could be used for future robustness research.

### 4.5 Which prompts are more robust? Word frequency

The word frequency results of these two datasets are presented in Appendix E. Our findings underscore that the resilience of a prompt is intricately tied to the contextual use of words, rather than the mere presence of certain terms. This complexity suggests that factors beyond word frequency, such as semantic coherence and syntactic structures, might be instrumental in determining robustness. This knowledge is valuable as it can influence future research on the robustness of LLMs, provide guidance for crafting more resistant prompts, and facilitate the creation of defensive mechanisms against adversarial prompt attacks. It is essential to emphasize that our observations are rooted in the current scope of models and datasets. Furthermore, the robustness or vulnerability of words remains deeply context-dependent.

Hence, direct determination of word robustness without considering the broader context may lead to oversimplified or inaccurate conclusions.

## 4.6 Countermeasures and defenses

We discuss potential countermeasures. **1) Input preprocessing:** One approach involves directly detecting and addressing potential adversaries, such as detecting typos, irrelevant sequences, and enhancing clarity and conciseness of prompts. **2) Incorporate low-quality data in pre-training:** Low-quality data can serve as potential adversaries, and explicitly including low-quality data during pre-training may develop a better understanding of diverse inputs and build resilience against adversaries. **3) Improved fine-tuning:** Fine-tuning could lead to improved robustness. As we demonstrated before, models such as T5 and UL2 exhibit greater robustness compared to ChatGPT, suggesting potential benefits of large-scale supervised fine-tuning. More details are in Appendix F.

## 5 Related work

**LLM Robustness Evaluation.** AdvGLUE (Wang et al., 2021) stands as a static dataset for evaluating adversarial robustness of *input samples*. DecodingTrust (Wang et al., 2023a) undertakes a comprehensive assessment of trustworthiness in GPT models, notably GPT-3.5 and GPT-4. The research delves into areas like toxicity, stereotype bias, adversarial challenges, and privacy. Specifically, they evaluate the robustness on standard datasets AdvGLUE(Wang et al., 2021) and AdvGLUE++ (Wang et al., 2023a). Specifically for adversarial robustness, DecodingTrust also focuses on evaluating the robustness of input samples instead of prompts and it still uses static datasets rather than an actionable benchmark suite. In contrast, PromptRobust is positioned as an open benchmark concentrating on adversarial *prompts* rather than samples (and it can be extended to include samples). Note that the prompts are general instructions to assist the in-context learning of LLMs to perform specific tasks, and they can be combined with many samples in certain tasks. Prompts are indispensable in human-LLMs interaction while input samples may not be needed, which means that prompts are versatile and it is essential to evaluate their robustness.

**Safety of LLMs.** We mimic the potential user prompts by creating adversarial prompts, but the main purpose is not to actually attack the model. This distinguishes our work from existing efforts in safety research of LLMs. Specifically, both SafetyPrompts (Sun et al., 2023) and prompt injection attacks (Perez & Ribeiro, 2022; Greshake et al., 2023; Liu et al., 2023b) are engineered to spotlight potentially harmful instructions that could steer LLMs into *delivering outputs misaligned with human values or perform unintended actions* such as data leakage and unauthorized access. Adversarial prompts are crafted to *mimic plausible mistakes an end-user might inadvertently make*. Our goal is to assess the extent to which these prompts, even if they slightly deviate from the norm, can skew LLM outcomes. These adversarial prompts retain their semantic integrity, ensuring they're virtually imperceptible for humans. The adversarial prompts are not designed to elicit harmful or misleading responses.

## 6 Conclusion and Limitation

The robustness of the prompts in LLMs is of paramount concern in security and human-computer interaction. In this paper, we thoroughly evaluated the robustness of LLMs to adversarial prompts using the proposed PromptRobust benchmark. The key is to leverage adversarial attack approaches to mimic potential perturbations such as typos, synonyms, and stylistic differences. We then performed extensive experiments and analysis on various tasks and models. While the results show that current LLMs are not robust enough to adversarial prompts, we further analyzed the reason behind it using attention visualization. Moreover, we analyzed the frequent words to provide guidance for both experts and non-experts in developing better prompt engineering tools. PromptRobust will be open-sourced to serve as a foundational tool for robust LLMs research.

There are several limitations. First, due to the substantial computation, we did not perform evaluations on the full datasets, but relied on sampling. Future research may evaluate on the entire datasets to gain more insights. Second, we cannot included all LLMs and datasets due to computation constraint. Including more in the future could provide a more diverse perspective. Third, we did not evaluate more advanced techniques of prompt engineering

such as chain-of-thought (CoT) (Wei et al., 2022) and tree-of-thought (ToT) (Yao et al., 2023). We believe more evaluations can be done on latest prompt engineering techniques. Fourth, we considered black-box prompt attacks, which can generate perturbations that can mimic naturally occurred errors. Optimized white-box prompt attacks may produce better adversarial prompts, which is an interesting future work.

# References

Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? FAccT 2021, pp. 610–623, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450383097.

Stella Biderman, USVSN Sai Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raf. Emergent and predictable memorization in large language models. *arXiv preprint arXiv:2304.11158*, 2023.

Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. Gpt-neox-20b: An open-source autoregressive language model, 2022. URL `https://arxiv.org/abs/2204.06745`.

Google Brain. A new open source flan 20b with ul2, 2023. URL `https://www.yitay.net/blog/flan-ul2-20b`.

Tom B Brown, Dandelion Mané, Aurko Roy, Martín Abadi, and Justin Gilmer. Adversarial patch. *arXiv preprint arXiv:1712.09665*, 2017.

Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations*, 2023.

Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Brian Strope, and Ray Kurzweil. Universal sentence encoder for English. In *EMNLP*, pp. 169–174, Brussels, Belgium, November 2018. Association for Computational Linguistics. doi: 10.18653/v1/ D18-2029. URL `https://aclanthology.org/D18-2029`.

Mauro Cettolo, Marcello Federico, Luisa Bentivogli, Jan Niehues, Sebastian Stüker, Katsuhito Sudoh, Koichiro Yoshino, and Christian Federmann. Overview of the IWSLT 2017 evaluation campaign. In *Proceedings of the 14th International Conference on Spoken Language Translation*, pp. 2–14, Tokyo, Japan, December 14-15 2017. International Workshop on Spoken Language Translation. URL `https://aclanthology.org/2017.iwslt-1.1`.

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. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL `https://lmsys.org/blog/2023-03-30-vicuna/`.

Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models, 2022. URL `https://arxiv.org/abs/2210.11416`.

Fred J Damerau. A technique for computer detection and correction of spelling errors. *Communications of the ACM*, 7(3):171–176, 1964.

Databricks. Hello dolly: Democratizing the magic of chatgpt with open models, 2023. URL `https://www.databricks.com/blog/2023/03/24/hello-dolly-democratizing-magic-chatgpt-open-models.html`.

Nolan Dey, Gurpreet Gosal, Zhiming, Chen, Hemant Khachane, William Marshall, Ribhu Pathria, Marvin Tom, and Joel Hestness. Cerebras-gpt: Open compute-optimal language models trained on the cerebras wafer-scale cluster, 2023.

William B. Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*, 2005. URL `https://aclanthology.org/I05-5002`.

Andreas Eisele and Yu Chen. MultiUN: A multilingual corpus from united nation documents. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, Valletta, Malta, May 2010. European Language Resources Association (ELRA). URL `http://www.lrec-conf.org/proceedings/lrec2010/pdf/686_Paper.pdf`.

Marzieh Fadaee, Arianna Bisazza, and Christof Monz. Data augmentation for low-resource neural machine translation. *arXiv preprint arXiv:1705.00440*, 2017.

J. Gao, J. Lanchantin, M. L. Soffa, and Y. Qi. Black-box generation of adversarial text sequences to evade deep learning classifiers. In *2018 IEEE Security and Privacy Workshops (SPW)*, pp. 50–56, May 2018. doi: 10.1109/SPW.2018.00016.

Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.

Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. More than you've asked for: A comprehensive analysis of novel prompt injection threats to application-integrated large language models. *arXiv preprint arXiv:2302.12173*, 2023.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021. URL `https://openreview.net/forum?id=d7KBjmI3GmQ`.

Jordan Hoffmann et al. Training compute-optimal large language models, 2022.

Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural computation*, 3(1):79–87, 1991.

Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. Is bert really robust? natural language attack on text classification and entailment. *arXiv preprint arXiv:1907.11932*, 2019.

Karen Kukich. Techniques for automatically correcting words in text. *ACM computing surveys (CSUR)*, 24(4):377–439, 1992.

Alexey Kurakin, Ian Goodfellow, Samy Bengio, et al. Adversarial examples in the physical world, 2016.

Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In *Thirteenth international conference on the principles of knowledge representation and reasoning*, 2012.

Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. TextBugger: Generating adversarial text against real-world applications. In *Proceedings 2019 Network and Distributed System Security Symposium*. Internet Society, 2019. doi: 10.14722/ndss.2019.23138. URL `https://doi.org/10.14722%2Fndss.2019.23138`.

Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. BERT-ATTACK: Adversarial attack against BERT using BERT. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6193–6202, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.500. URL `https://aclanthology.org/2020.emnlp-main.500`.

Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. Evaluating the logical reasoning ability of chatgpt and gpt-4. *arXiv preprint arXiv:2304.03439*, 2023a.

Yi Liu, Gelei Deng, Yuekang Li, Kailong Wang, Tianwei Zhang, Yepang Liu, Haoyu Wang, Yan Zheng, and Yang Liu. Prompt injection attack against llm-integrated applications. *arXiv preprint arXiv:2306.05499*, 2023b.

Inbal Magar and Roy Schwartz. Data contamination: From memorization to exploitation. *arXiv preprint arXiv:2203.08242*, 2022.

Pankaj Malhotra, Lovekesh Vig, Gautam Shroff, Puneet Agarwal, et al. Long short term memory networks for anomaly detection in time series. In *Esann*, volume 2015, pp. 89, 2015.

Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *CVPR*, pp. 1765–1773, 2017.

John Morris, Eli Lifland, Jin Yong Yoo, Jake Grigsby, Di Jin, and Yanjun Qi. Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in nlp. In *EMNLP*, pp. 119–126, 2020.

Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. Stress test evaluation for natural language inference. In *ACL*, pp. 2340–2353, Santa Fe, New Mexico, USA, August 2018. Association for Computational Linguistics. URL `https://aclanthology.org/C18-1198`.

Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial nli: A new benchmark for natural language understanding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4885–4901, 2020.

OpenAI. `https://chat.openai.com.chat`, 2023a.

OpenAI. Gpt-4 technical report, 2023b.

Tianyu Pang, Kun Xu, Chao Du, Ning Chen, and Jun Zhu. Improving adversarial robustness via promoting ensemble diversity. In *International Conference on Machine Learning*, pp. 4970–4979. PMLR, 2019.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *ACL*, pp. 311–318, July 2002. doi: 10.3115/1073083.1073135. URL `https://aclanthology.org/P02-1040`.

Fábio Perez and Ian Ribeiro. Ignore previous prompt: Attack techniques for language models. *arXiv preprint arXiv:2211.09527*, 2022.

Massimo Quadrana, Alexandros Karatzoglou, Balázs Hidasi, and Paolo Cremonesi. Personalizing session-based recommendations with hierarchical recurrent neural networks. In *proceedings of the Eleventh ACM Conference on Recommender Systems*, pp. 130–137, 2017.

Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don't know: Unanswerable questions for SQuAD. In *ACL*, pp. 784–789, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-2124. URL `https://aclanthology.org/P18-2124`.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Semantically equivalent adversarial rules for debugging nlp models. In *Proceedings of the 56th annual meeting of the association for computational linguistics (volume 1: long papers)*, pp. 856–865, 2018.

Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. Beyond accuracy: Behavioral testing of NLP models with CheckList. In *ACL*, pp. 4902–4912, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.442. URL https://aclanthology.org/2020.acl-main.442.

David Saxton, Edward Grefenstette, Felix Hill, and Pushmeet Kohli. Analysing mathematical reasoning abilities of neural models. In *ICLR*, 2019. URL https://openreview.net/forum?id=H1gR5iR5FX.

Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *EMNLP*, pp. 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics. URL https://www.aclweb.org/anthology/D13-1170.

Hao Sun, Zhexin Zhang, Jiawen Deng, Jiale Cheng, and Minlie Huang. Safety assessment of chinese large language models. *arXiv preprint arXiv:2304.10436*, 2023.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023a.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. Universal adversarial triggers for attacking and analyzing nlp. *arXiv preprint arXiv:1908.07125*, 2019.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. 2019. In the Proceedings of ICLR.

Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. https://github.com/kingoflolz/mesh-transformer-jax, May 2021.

Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and Bo Li. Adversarial glue: A multi-task benchmark for robustness evaluation of language models. *arXiv preprint arXiv:2111.02840*, 2021.

Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. *arXiv preprint arXiv:2306.11698*, 2023a.

Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen, Runkai Zheng, Yidong Wang, Linyi Yang, Haojun Huang, Wei Ye, Xiubo Geng, et al. On the robustness of chatgpt: An adversarial and out-of-distribution perspective. In *International conference on learning representations (ICLR) workshop on Trustworthy and Reliable Large-Scale Machine Learning Models*, 2023b.

Zhiguo Wang, Wael Hamza, and Radu Florian. Bilateral multi-perspective matching for natural language sentences, 2017.

Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. Neural network acceptability judgments. *arXiv preprint arXiv:1805.12471*, 2018.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022.

Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for sentence understanding through inference. In *NAACL HLT*, pp. 1112–1122. Association for Computational Linguistics, 2018. URL `http://aclweb.org/anthology/N18-1101`.

Thomas Wolf et al. Huggingface's transformers: State-of-the-art natural language processing, 2020.

Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International Conference on Machine Learning*, pp. 23965–23998. PMLR, 2022a.

Mitchell Wortsman, Gabriel Ilharco, Jong Wook Kim, Mike Li, Simon Kornblith, Rebecca Roelofs, Raphael Gontijo Lopes, Hannaneh Hajishirzi, Ali Farhadi, Hongseok Namkoong, et al. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7959–7971, 2022b.

Linyi Yang, Shuibai Zhang, Libo Qin, Yafu Li, Yidong Wang, Hanmeng Liu, Jindong Wang, Xing Xie, and Yue Zhang. Glue-x: Evaluating natural language understanding models from an out-of-distribution generalization perspective. In *ACL 2023 Findings*, 2023.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of Thoughts: Deliberate problem solving with large language models, 2023.

Yuan Zang, Fanchao Qi, Chenghao Yang, Zhiyuan Liu, Meng Zhang, Qun Liu, and Maosong Sun. Word-level textual adversarial attacking as combinatorial optimization. In *Proceedings of ACL*, 2020.

Ziqi Zhang, Yuanchun Li, Jindong Wang, Bingyan Liu, Ding Li, Yao Guo, Xiangqun Chen, and Yunxin Liu. Remos: reducing defect inheritance in transfer learning via relevant model slicing. In *Proceedings of the 44th International Conference on Software Engineering*, pp. 1856–1868, 2022.

Terry Yue Zhuo, Zhuang Li, Yujin Huang, Yuan-Fang Li, Weiqing Wang, Gholamreza Haffari, and Fatemeh Shiri. On robustness of prompt-based semantic parsing with large pre-trained language model: An empirical study on codex. *arXiv preprint arXiv:2301.12868*, 2023.

## A Details on Adversarial Prompts

### A.1 Examples of prompts

Table 7 shows the details of the task and role-oriented prompts in zero-shot and few-shot settings. Table 8 shows examples of 7 adversarial attacks on prompts.

Table 7: Examples of 4 types of prompts.

| | | |
|---|---|---|
| Zero-shot | Task-oriented | Evaluate the sentiment of the given text and classify it as 'positive' or 'negative': |
| | Role-oriented | In the role of a sentiment analysis tool, respond with 'positive' or 'negative' to classify this statement: |
| Few-shot | Task-oriented | Analyze the tone of this statement and respond with either 'positive' or 'negative'. Here are three examples. Sentence: hide new secretions from the parental units. Answer: negative. Sentence: contains no wit , only labored gags. Answer: negative. Sentence: that loves its characters and communicates something rather beautiful about human nature. Answer: positive. |
| | Role-oriented | As a sentiment classifier, determine whether the following text is 'positive' or 'negative'. Here are three examples. Sentence: hide new secretions from the parental units. Answer: negative. Sentence: contains no wit , only labored gags. Answer: negative. Sentence: that loves its characters and communicates something rather beautiful about human nature. Answer: positive. |

Table 8: Example of adversarial prompts generated by 7 prompt attacks to mimic possible prompts. The characters and words marked with red are generated by prompt attacks.

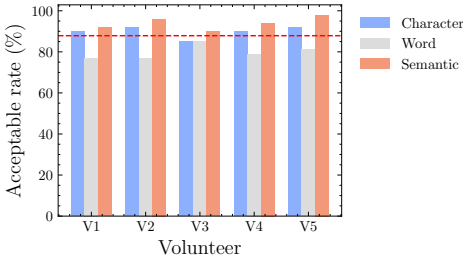| Scenario | Prompt |
|---|---|
| Clean | `As a mathematics instructor, calculate the answer to the following problem related to {}:` |
| TextBugger | `As a mathematics instructorr, calculate the answers to the following problem related to {}:` |
| DeepWordBug | `As a mathematics iestructor, calculate the answex to the following problem related to {}:` |
| TextFooler | `As a mathematics prof, calculate the address to the following problem related to {}:` |
| BertAttack | `As a mathematics instructor, calculate the sum to the following problem related to {}:` |
| CheckList | `As a mathematics instructor, calculate the answer to the following problem related to KjPJJ2a7RB {}:` |
| StressTest | `As a mathematics instructor, calculate the answer to the following problem related to and false is not true {}:` |
| Semantic | `Compute the result of {}.` |



Figure 2: Results of human study on semantic preserving of the adversarial prompts. The dotted red line ($> 85\%$) is the average of all volunteers on all attacks.

## A.2 Semantic preservation of adversarial prompts

In our endeavor to validate the efficacy of our adversarial prompt generation, we engaged in a human-centric study. We enlisted the expertise of 5 independent evaluators with proficiency in the domain to critically assess the semantic congruence between the original and the generated adversarial prompts. For the study, we randomly sampled a set of 100 adversarial prompts along with their respective original prompts. The evaluators were tasked with determining if each adversarial prompt was semantically equivalent to its original counterpart, ensuring the absence of semantic drift. Such an evaluation provides insights into the reliability and robustness of our adversarial prompts.

To address the challenges associated with word-level attacks, we have diligently fine-tuned the hyperparameters of each attack approach, thus striving to maintain semantic continuity. Then, we conduct a human study to recruit five volunteers to judge if the generated adversarial prompts can preserve semantics. The evaluators were presented with the original prompt $P$ juxtaposed with its adversarial version $\bar{P}$, and were tasked with determining their semantic congruence. Sentence-level attacks are excluded in this study since they do not change the original prompts, but only to add extra perturbations in the end. The results in Figure 2 demonstrate that these adversarial prompts generated by character-level, word-level and semantic-level attacks are **at least** 85% **acceptable by humans**, indicating that our attack is realistic and meaningful.

Here, we present a few exemplars of both acceptable and unacceptable adversarial prompts as discerned by the evaluators:

- **SST-2** Original prompt: `Read the provided excerpt and choose between 'positive' and 'negative' to describe its sentiment:`
  Attacked prompt: `read the standard excerpt and pause between 'positive' and 'negative' to describe its thoughtful:`
  Unacceptable.
- **CoLA** Original prompt: `Functioning as a grammar evaluation tool, analyze the given sentence and decide if it is grammatically correct, responding with 'acceptable' or 'unacceptable':`

Attacked prompt: `Functioning as a grammar evaluation Lool, analyze the Tgiven sentence and deciRe if it is grammXatically correct, respoOnding with 'acceptable' or 'unacceptable':`
<span style="color:green">Acceptable.</span>

- **MMLU** Original prompt: `As a knowledgeable individual in , provide your response to the multiple-choice question by choosing 'A', 'B', 'C', or 'D'.`
Attacked prompt: `As a knowledgeable everybody in , prescribes your rebuttal to the multiple-choice item by chooses 'A', 'B', 'C', or 'D'.`
<span style="color:red">Unacceptable.</span>

- **Math** Original prompt: `Resolve the following mathematical question associated with:`
Attacked prompt: `Resolve the following mathematical answer along with :`
<span style="color:green">Acceptable.</span>

## A.3   Prompt attack

The majority of our prompt attacks have been developed by adapting and revising strategies from TextAttack[4] (Morris et al., 2020). For the detailed settings of each attack, please refer to our code.

**Character Level:** Techniques such as TextBugger and DeepWordBug manipulate text at the character level by introducing typos or errors within words through insertions, deletions, replacements, and replications. These methods capitalize on the model's vulnerability to minor perturbations in individual characters, frequently resulting in misclassification or erroneous interpretations.

We primarily adopt the settings from TextAttack for TextBugger and DeepWordBug, such as the repeat constraint which prohibits modifying words that have already been altered. Additionally, For TextBugger, TextAttack enforces a constraint on the overall similarity between the sentence encodings of clean and adversarial prompts, utilizing the Universal Sentence Encoder (Cer et al., 2018) to generate text embeddings. In our study, we set this minimum similarity threshold to 0.8. For DeepWordBug, TextAttack set constraint on edit distance (Levenshtein Distance) as 30.

**Word Level:** In this study, we employ BertAttack and TextFooler for word-level attacks. These approaches focus on replacing words within the text with synonyms or contextually similar words. By making ostensibly minor alterations to the input text, these attacks can deceive large language models into producing incorrect outputs or substantially modifying their predictions. We meticulously fine-tune the hyperparameters of BertAttack and TextFooler to obtain more appropriate synonyms.

For TextFooler, we set the minimum embedding cosine similarity between word and its synonyms as 0.6, and the minimum Universal Sentence Encoder similarity is 0.84. For BertAttack, the minimum Universal Sentence Encoder similarity is 0.8.

**Sentence Level:** StressTest and CheckList serve as examples of sentence-level attacks, wherein adversaries attempt to distract the model by adding irrelevant or extraneous sentences to the input text. By incorporating misleading information into the text, these methods can potentially cause the model to lose focus on the primary context, leading to inaccurate results. For the StressTest attack, we adopt similar settings to those in (Wang et al., 2019), appending "`and true is true, `" "`and false is not true, `" or "`and true is true `" for five times to the end of a prompt. For the CheckList attack, we generate 50 random sequences consisting of alphabets and digits, each with a length of 10, and append this random sequences into the end of a prompt.

**Semantic Level:** At the human level, adversaries can construct prompts using various languages, such as Chinese, French, Arabic, Spanish, Japanese, and Korean, subsequently translating these prompts into English. By exploiting the nuances and idiosyncrasies of

---

[4] `https://github.com/QData/TextAttack`

different languages during translation, it can introduce subtle ambiguities, grammatical errors, or inconsistencies in the input prompt. This poses a formidable challenge for NLP models in generating accurate and coherent responses.

For each language, we first construct 10 prompts based on a English prompt by GPT4 (OpenAI, 2023b), then translate it back to English by Google Translator.

## A.4   Our attack vs. existing textual adversarial attacks

Prompt attacks and textual adversarial attacks (Li et al., 2019; Jin et al., 2019; Gao et al., 2018; Li et al., 2020; Ribeiro et al., 2020; Naik et al., 2018; Zang et al., 2020) are both rooted in similar foundational algorithms, but differ in critical ways:

- Target of attack: Prompt attacks target the *instruction (prompts)* for LLMs while vanilla adversarial attacks focus on the *samples*. In numerous tasks, the data might be optional, while prompts remain indispensable. For example, "Write a story about a fox." and "Give me some investigation suggestions." are all prompts with *no* samples. This makes our prompt attack more general.
- Universality of adversarial prompts: An adversarial prompt, represented as $\bar{P}$, works as a common threat for all samples related to a specific task. For example, if $P$ is designed to instruct LLMs to solve math problems, then $\bar{P}$ can be used for *many* different math problems and datasets. This ability is significantly different from current NLP adversarial benchmarks.

In essence, prompt attacks seek to delve into the universality (Moosavi-Dezfooli et al., 2017; Wallace et al., 2019) of adversarial prompts. We argue this offers an innovative lens to assess the robustness of language models, complementing insights from existing benchmarks like AdvGLUE (Wang et al., 2021), and AdvGLUE++(Wang et al., 2023a).

# B   Models, Datasets, and Environments

## B.1   Models

Here, we list the brief introduction of each LLM in our experiments. For more details about Vicuna, please refer to its GitHub repository[5]. For the other LLMs, please refer to Huggingface transformer repository (Wolf et al., 2020).

- **Flan-T5-large (Chung et al., 2022)**: Flan-T5-large is a derivative of the Text-to-Text Transfer Transformer (T5) model, developed by Google.
- **Dolly-6B (Databricks, 2023)**: The Dolly-v1-6b model is a 6-billion parameter causal language model developed by Databricks. It originates from EleutherAI's GPT-J (Wang & Komatsuzaki, 2021) and has been fine-tuned on the Stanford Alpaca (Taori et al., 2023) corpus, which comprises roughly 52K question/answer pairs.
- **Vicuna-13B (Chiang et al., 2023)**: Vicuna-13B, fine-tuned from the LLaMA-13B base model, was developed using approximately 70K user-shared conversations collected from ShareGPT.com via public APIs.
- **Cerebras-13B (Dey et al., 2023)**: Cerebras-13B is based on the GPT-3 style architecture. All models in the Cerebras-GPT series have been trained according to Chinchilla scaling laws (Hoffmann et al., 2022), which optimize compute efficiency by maintaining a ratio of 20 tokens per model parameter.
- **Llama2-13B (Touvron et al., 2023a)**: The Llama2 model, developed by the FAIR team at Meta AI, is an autoregressive language model that employs the transformer architecture.
- **GPT-NEOX-20B (Black et al., 2022)**: GPT-NEOX-20B is a large-scale implementation of GPT-based models, with NEOX-20B specifically referring to a variant of this series comprising 20 billion parameters.

---

[5]`https://github.com/lm-sys/FastChat`

- **Flan-UL2 (Brain, 2023)**: Flan-UL2 is an encoder decoder model based on the T5 architecture. It uses the same configuration as the UL2 model. It was fine-tuned using the "Flan" prompt tuning and dataset collection.

- **ChatGPT (OpenAI, 2023a) and GPT-4 (OpenAI, 2023b)**: Developed by OpenAI, ChatGPT is a large language model trained to generate human-like text based on the prompt it's given. It uses the GPT-3 architecture and has been fine-tuned for more interactive and conversational tasks. GPT-4 is by far the best-performing LLMs.

## B.2   Tasks and datasets

We adopt the following public datasets for evaluation and PromptRobust can easily take as inputs other datasets.

- **GLUE** The GLUE dataset (General Language Understanding Evaluation) (Wang et al., 2019) is a collection of resources designed to assess and benchmark the performance of natural language processing (NLP) models across various language understanding tasks. In this study, we selected 8 tasks, including Sentiment Analysis (SST-2 (Socher et al., 2013)), Grammar Correctness (CoLA (Warstadt et al., 2018)), Duplicate Sentence Detection (QQP (Wang et al., 2017), MRPC (Dolan & Brockett, 2005)), and Natural Language Inference (MNLI (Williams et al., 2018), QNLI (Wang et al., 2019), RTE (Wang et al., 2019), and WNLI (Levesque et al., 2012)).

- **MMLU (Hendrycks et al., 2021)** To evaluate the extensive world knowledge and problem-solving abilities of large language models, the MMLU dataset encompasses 57 tasks consisting of multiple-choice questions from diverse domains, such as mathematics, history, computer science, law, and more. This dataset serves as a massive multitask test.

- **SQuAD V2 (Rajpurkar et al., 2018)** SQuAD v2 is a widely used dataset for training and evaluating natural language processing models in the domain of machine reading comprehension. SQuAD v2 enhances the original SQuAD dataset (SQuAD v1) by introducing unanswerable questions, increasing the challenge for models. For each question, the model must either (1) identify the correct answer span within the passage (if the question is answerable) or (2) predict that the question is unanswerable (if there is no answer span within the passage).

- **UN Multi (Eisele & Chen, 2010)** The Multi UN dataset is a large parallel corpus of text gathered from official United Nations documents. It comprises texts in six official languages of the United Nations: Arabic, Chinese, English, French, Russian, and Spanish. The Multi UN dataset primarily contains formal texts, which may limit its applicability to more informal language domains or conversational applications.

- **IWSLT 2017 (Cettolo et al., 2017)** The IWSLT 2017 dataset (International Workshop on Spoken Language Translation 2017) is a collection of multilingual, multi-domain parallel text data specifically designed for evaluating spoken language translation systems. The translation tasks include data from the TED Talks Open Translation Project, featuring parallel text data for multiple language pairs such as English-German, English-French, English-Chinese, and English-Czech. The dataset consists of both spoken language transcriptions and their corresponding translations.

- **Math (Saxton et al., 2019)** DeepMind Mathematics Dataset is a collection of math problems aimed at evaluating the mathematical reasoning abilities of artificial intelligence models. The dataset challenges AI models to solve a diverse range of mathematical problems, spanning from algebra to calculus, and tests their ability to comprehend and reason via complex mathematical concepts.

## B.3   Environments

To reproduce the computational environment used in this study, an environment file, `environment.yml`, is provided in our repository. This YAML file lists all the dependencies and their specific versions used in the study. Users can create an identical Conda environment using the command `conda env create -f environment.yml`. The computational

experiments were conducted on machines equipped with NVIDIA Tesla V100 GPUs (16GB of GPU memory each).

## C   Details on Experiments and Results

### C.1   Details on dataset sampling

Note that we cannot run evaluations on all samples due to significantly extensive computing requirements. Instead, we turn to sampling. Specifically, for the GLUE datasets, we sample 1,000 instances when the validation set exceeds this size; otherwise, we utilize the entire validation set. With respect to ChatGPT and GPT4, we adopt a smaller sample size of 200 instances for computational efficiency. For the MMLU dataset, we select 10 instances for each of the 57 tasks if the validation set exceeds this size; if not, the entire validation set is used. For the SQUAD V2 dataset, we randomly select 200 validation instances. Regarding the translation datasets UN Multi and IWSLT 2017, we focus on three languages—English, French, and German, which are primarily supported by T5-large and UL2. We select a total of 100 validation instances, evenly distributed among all possible translation pairs, e.g., English to French. For the Math dataset, we select 20 types of math problems, choosing either 5 or 10 instances per type, resulting in a total of 160 instances. This sampling strategy ensures the formation of a manageable and representative evaluation set for each dataset, thereby enabling an effective assessment of the performance and robustness of LLMs across various tasks and domains.

### C.2   Results of clean prompts on all LLMs

Table 9 showcases the performance of different models across various datasets when using clean prompts. Certain LLMs, including Dolly, Cerebras, and NEXO, encounter difficulties with some datasets. For instance, Dolly's accuracy for the QQP dataset is merely 0.53%, a stark contrast to T5's accuracy of 86.67%. Consequently, we focus our attack study on models that demonstrate superior performance, namely T5, Vicuna, Llama2, UL2, ChatGPT, and GPT4.

Table 9: The Average performance and standard deviations of different models on different datasets.

| Dataset | T5 | Dolly | Vicuna | Cerebras | Llama2 | NEOX | UL2 | ChatGPT |
|---|---|---|---|---|---|---|---|---|
| SST-2 | $94.79_{\pm0.49}$ | $47.80_{\pm9.30}$ | $21.12_{\pm15.40}$ | $21.33_{\pm23.02}$ | $90.25_{\pm2.23}$ | $21.49_{\pm13.35}$ | $95.92_{\pm1.03}$ | $92.91_{\pm3.32}$ |
| CoLA | $76.11_{\pm1.28}$ | $4.92_{\pm9.04}$ | $35.28_{\pm20.12}$ | $18.18_{\pm23.82}$ | $74.53_{\pm1.87}$ | $7.96_{\pm14.23}$ | $86.07_{\pm0.36}$ | $78.91_{\pm1.75}$ |
| QQP | $86.67_{\pm1.05}$ | $0.53_{\pm1.66}$ | $24.74_{\pm10.03}$ | $0.00_{\pm0.00}$ | $23.23_{\pm6.97}$ | $0.00_{\pm0.02}$ | $88.25_{\pm0.54}$ | $81.49_{\pm1.47}$ |
| MRPC | $80.75_{\pm1.73}$ | $0.17_{\pm0.30}$ | $50.15_{\pm19.65}$ | $0.01_{\pm0.05}$ | $49.15_{\pm4.56}$ | $0.01_{\pm0.05}$ | $86.03_{\pm1.41}$ | $72.71_{\pm2.82}$ |
| MNLI | $81.39_{\pm4.7}$ | $0.78_{\pm0.88}$ | $12.90_{\pm8.21}$ | $0.87_{\pm1.16}$ | $57.30_{\pm1.53}$ | $0.00_{\pm0.00}$ | $83.50_{\pm4.79}$ | $76.71_{\pm2.44}$ |
| QNLI | $85.12_{\pm5.57}$ | $0.05_{\pm0.07}$ | $27.76_{\pm10.04}$ | $0.00_{\pm0.00}$ | $14.90_{\pm8.48}$ | $4.22_{\pm5.46}$ | $93.68_{\pm0.41}$ | $77.53_{\pm7.48}$ |
| RTE | $84.24_{\pm1.16}$ | $0.19_{\pm0.77}$ | $29.51_{\pm15.12}$ | $0.00_{\pm0.00}$ | $47.67_{\pm1.92}$ | $3.16_{\pm4.40}$ | $93.26_{\pm0.51}$ | $80.73_{\pm3.24}$ |
| WNLI | $62.34_{\pm3.31}$ | $0.00_{\pm0.00}$ | $22.57_{\pm15.96}$ | $0.00_{\pm0.00}$ | $41.08_{\pm1.71}$ | $3.62_{\pm5.10}$ | $77.53_{\pm1.4}$ | $61.07_{\pm6.22}$ |
| MMLU | $45.25_{\pm0.83}$ | - | $15.31_{\pm7.41}$ | - | $36.05_{\pm7.76}$ | - | $53.04_{\pm0.67}$ | $63.33_{\pm2.56}$ |
| SQuAD V2 | $87.32_{\pm0.43}$ | - | - | - | - | - | $89.78_{\pm0.71}$ | $68.35_{\pm4.36}$ |
| IWSLT | $0.18_{\pm0.04}$ | - | - | - | - | - | $0.21_{\pm0.04}$ | $0.23_{\pm0.01}$ |
| UN Multi | $0.29_{\pm0.02}$ | - | - | - | - | - | $0.33_{\pm0.02}$ | $0.34_{\pm0.01}$ |
| Math | $14.22_{\pm3.25}$ | - | - | - | - | - | $14.81_{\pm1.35}$ | $13.14_{\pm8.48}$ |

### C.3   Analysis on model size and fine-tuning

As shown in Table 10 and 3, there seems to be no clear correlation between model robustness and size, for example, despite being the smallest, T5-large demonstrates robustness on par with larger models such as ChatGPT on our evaluated datasets.

The observed differences in model robustness might stem from two aspects: 1) the specific fine-tuning techniques employed. For example, both UL2 and T5-large, fine-tuned on large datasets, and ChatGPT, fine-tuned via RLHF (Christiano et al., 2017), exhibit better robustness than Vicuna. These findings encourage further investigation of fine-tuning strategies to enhance robustness. 2) the memorization of training data. Recent studies suggest that the remarkable performance of some LLMs might be rooted in their ability to memorize training data (Bender et al., 2021; Magar & Schwartz, 2022; Carlini et al.,
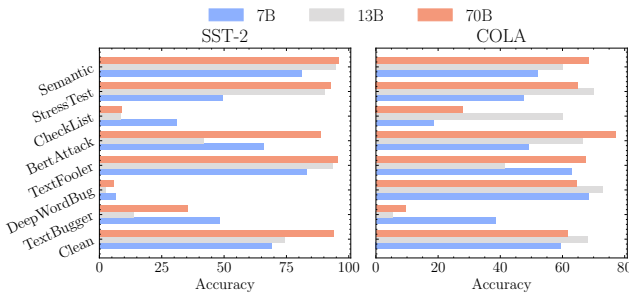
Table 10: The APDR and standard deviations of different attacks on different models.

| Model | Character-level | | Word-level | | Sentence-level | | Semantic-level |
|---|---|---|---|---|---|---|---|
| | TextBugger | DeepWordBug | TextFooler | BertAttack | CheckList | StressTest | Semantic |
| T5-large | $0.09_{\pm0.10}$ | $0.13_{\pm0.18}$ | $0.20_{\pm0.24}$ | $0.21_{\pm0.24}$ | $0.04_{\pm0.08}$ | $0.18_{\pm0.24}$ | $0.10_{\pm0.09}$ |
| Vicuna | $0.81_{\pm0.25}$ | $0.69_{\pm0.30}$ | $0.80_{\pm0.26}$ | $0.84_{\pm0.23}$ | $0.64_{\pm0.27}$ | $0.29_{\pm0.40}$ | $0.74_{\pm0.25}$ |
| Llama2 | $0.67_{\pm0.36}$ | $0.41_{\pm0.34}$ | $0.68_{\pm0.36}$ | $0.74_{\pm0.33}$ | $0.34_{\pm0.33}$ | $0.20_{\pm0.30}$ | $0.66_{\pm0.35}$ |
| UL2 | $0.04_{\pm0.06}$ | $0.03_{\pm0.04}$ | $0.14_{\pm0.20}$ | $0.16_{\pm0.22}$ | $0.04_{\pm0.07}$ | $0.06_{\pm0.09}$ | $0.06_{\pm0.08}$ |
| ChatGPT | $0.14_{\pm0.20}$ | $0.08_{\pm0.13}$ | $0.32_{\pm0.35}$ | $0.34_{\pm0.34}$ | $0.07_{\pm0.13}$ | $0.06_{\pm0.12}$ | $0.26_{\pm0.22}$ |
| GPT-4 | $0.03_{\pm0.10}$ | $0.02_{\pm0.08}$ | $0.18_{\pm0.19}$ | $0.27_{\pm0.40}$ | $-0.02_{\pm0.09}$ | $0.03_{\pm0.15}$ | $0.03_{\pm0.16}$ |
| Avg | $0.21_{\pm0.30}$ | $0.17_{\pm0.26}$ | $0.31_{\pm0.33}$ | $0.33_{\pm0.34}$ | $0.12_{\pm0.23}$ | $0.11_{\pm0.23}$ | $0.22_{\pm0.26}$ |

2023; Biderman et al., 2023), rather than in generalization. Hence, even when confronted with adversarial prompts, models might leverage this memorization to produce accurate responses.

In this section, we conduct experiments to analyze the effects of different model sizes and fine-tuning on adversarial prompts. Particularly, we leverage the open-source Llama2 (Touvron et al., 2023b) series due to their support on different sizes and their corresponding fine-tuned versions. The chat versions of Llama2 (Llama2-chat) are fine-tuned on human instructions datasets to better follow the instructions and support multi-turn conversations, while the original version can only be used for inference.

**Robustness of different model sizes** We analyzed three models from the open-source Llama2 series (Touvron et al., 2023b): Llama2-7B-chat, Llama2-13B-chat, and Llama2-70B-chat. These were chosen due to their distinct sizes, further, their fine-tuning datasets and methods are the same. Our results, depicted in Figure 3(a), reveal that, in a non-adversarial setting, larger models like the 70B model typically surpass the performance of their smaller counterparts. Yet, when subjected to adversarial attacks, the performance dynamics change: at times, smaller models outshine the larger ones. The reasons for these abnormal behaviors could trigger interests for future research.



(a) Analysis of model size.



(b) Analysis of fine-tuning.

Figure 3: (a) Accuracy of Llama2 models (7B-chat, 13B-chat, 70B-chat) on SST2 and CoLA datasets. (b) Accuracy of Llama2 models with fine-tuning and w/o fine-tuning (vanilla) on SST-2 dataset.

**Robustness of fine-tuning** To delve into the intricacies of fine-tuning, we compared the performances of Llama2-7B and Llama2-7B-chat on the SST2 and COLA tasks. Our analysis, as visualized in Figure 3(b), underscores a consistent trend: models fine-tuned using human-instruction datasets fare better against adversarial onslaughts than models that are not fine-tuned. This observation implies that fine-tuning could be further utilized as the countermeasures for adversarial inputs.

Table 11: Accuracy (%) of GPT-4 on clean and adversarial prompts and samples on the AdvGLUE Wang et al. (2021) dataset, i.e., attacking both prompts and samples.

| Attack | SST-2 | QQP | MNLI | QNLI | RTE | AVG |
|---|---|---|---|---|---|---|
| Clean prompts & clean samples | 96.10 | 78.23 | 81.05 | 64.50 | 87.54 | 81.48 |
| Clean prompts & AdvGLUE | 63.51 | 70.51 | 63.64 | 62.84 | 74.07 | 66.91 |
| TextBugger | 58.78 | 44.87 | 47.93 | 60.81 | 76.54 | 57.79 |
| DeepWordBug | 66.22 | 61.54 | 59.50 | 61.49 | 72.84 | 64.32 |
| TextFooler | 2.03 | 1.28 | 46.28 | 4.05 | 0.00 | 10.73 |
| BertAttack | 0.00 | 0.00 | 27.27 | 24.32 | 71.60 | 24.64 |
| Checklist | 69.59 | 66.67 | 57.85 | 56.08 | 72.84 | 64.61 |
| StressTest | 50.68 | 56.41 | 59.50 | 59.46 | 76.54 | 60.52 |
| Semantic | 0.00 | 38.46 | 48.76 | 65.54 | 41.98 | 38.95 |

## C.4 Results excluding non-semantic preserving adversarial prompts

Table 12 presents the attack results after excluding adversarial prompts that do not preserve semantics. It can be observed that the APDR is still considerably high for each dataset.

## C.5 Attacking both prompts and samples

The primary focus of this work is to evaluate the robustness of prompts rather than input samples since the samples can be omitted in certain situations, as discussed in Sec. 2.2. In this section, we explore the possibility of attacking both prompts and samples, i.e., evaluating the performance of LLMs in adversarial prompts and samples. Note that since the generation of adversarial examples is expensive and time-consuming, we leverage an existing adversarial dataset called AdvGLUE (Wang et al., 2021), which contains adversarial examples from GLUE (Wang et al., 2019) and it consists of five same tasks as GLUE: SST-2, QQP, MNLI, QNLI, and RTE. Then, we leverage the adversarial prompts and the AdvGLUE dataset (Wang et al., 2021) to evaluate the performance when attacking both prompts and samples.

Table **??** shows the accuracy results using both clean and adversarial prompts on AdvGLUE and clean dataset, respectively. The results demonstrate that on average, all attacking approaches are effective since the accuracy is dramatically declined in face of adversarial prompts. Similar to Table 2, word-level attacks (TextFooler and BertAttack) are the most effective with more than 49% of accuracy drop. Moreover, surprising results emerge for Checklist attack since the performance can sometimes be improved (e.g., 69.59% on SST-2 vs. the original 63.51%). This is also consistent with our previous observation in Sec. D.3. The results in this section show that attacking both the prompts and samples can further reduce the performance of LLMs. However, certain attacks can even enhance the performance, which is left for future research.

# D    Effectiveness Analysis of LLM Response and Attention Visualization

## D.1 Erroneous response analysis

We first analyze the erroneous response analysis produced by adversarial prompts. The results suggest that adversarial prompts impact LLMs' performance by inducing misclassification errors and hindering their ability to generate meaningful responses.

Table 12: The APDR and standard deviations of different attacks on different datasets by *excluding* the ones human annotators do not think acceptable.

| Dataset | Character-level | | Word-level | | Sentence-level | | Semantic-level |
|---|---|---|---|---|---|---|---|
| | TextBugger | DeepWordBug | TextFooler | BertAttack | CheckList | StressTest | Semantic |
| SST-2 | 0.26±0.39 | 0.21±0.36 | 0.33±0.35 | 0.30±0.39 | 0.27±0.39 | 0.17±0.34 | 0.28±0.36 |
| CoLA | 0.37±0.39 | 0.29±0.36 | 0.40±0.33 | 0.42±0.31 | 0.25±0.32 | 0.21±0.28 | 0.27±0.35 |
| QQP | 0.20±0.32 | 0.18±0.27 | 0.26±0.31 | 0.29±0.33 | 0.13±0.25 | -0.00±0.21 | 0.30±0.36 |
| MRPC | 0.24±0.33 | 0.21±0.30 | 0.27±0.30 | 0.31±0.29 | 0.13±0.27 | 0.20±0.30 | 0.28±0.36 |
| MNLI | 0.26±0.37 | 0.18±0.31 | 0.27±0.36 | 0.34±0.32 | 0.16±0.26 | 0.11±0.27 | 0.11±0.04 |
| QNLI | 0.36±0.39 | 0.41±0.36 | 0.47±0.33 | 0.45±0.30 | 0.22±0.37 | 0.18±0.26 | 0.35±0.33 |
| RTE | 0.24±0.37 | 0.22±0.36 | 0.26±0.34 | 0.28±0.35 | 0.19±0.32 | 0.18±0.25 | 0.28±0.33 |
| WNLI | 0.28±0.36 | 0.26±0.35 | 0.27±0.31 | 0.28±0.29 | 0.19±0.30 | 0.19±0.26 | 0.36±0.32 |
| MMLU | 0.18±0.22 | 0.11±0.15 | 0.19±0.16 | 0.31±0.20 | 0.14±0.20 | 0.03±0.16 | 0.17±0.17 |
| SQuAD V2 | 0.09±0.17 | 0.05±0.08 | 0.23±0.25 | 0.30±0.29 | 0.02±0.03 | 0.02±0.04 | 0.07±0.09 |
| IWSLT | 0.09±0.14 | 0.11±0.12 | 0.26±0.25 | 0.12±0.16 | 0.10±0.10 | 0.17±0.19 | 0.18±0.14 |
| UN Multi | 0.06±0.08 | 0.08±0.12 | 0.17±0.19 | 0.10±0.13 | 0.06±0.07 | 0.09±0.11 | 0.15±0.18 |
| Math | 0.19±0.17 | 0.15±0.13 | 0.45±0.32 | 0.39±0.27 | 0.16±0.11 | 0.13±0.08 | 0.23±0.13 |
| Avg | 0.23±0.33 | 0.20±0.30 | 0.29±0.31 | 0.31±0.30 | 0.16±0.27 | 0.13±0.25 | 0.24±0.29 |

Table 13: Attention visualization of samples that are *correctly classified by adv. prompt but misclassified by clean prompt.* Notations and colors follow Table 5.



- **Induced Misclassification:** As exemplified by BertAttack, CheckList, and Translation attacks, adversarial prompts can lead the model to erroneous classifications. For instance, the sentiment prediction may shift from positive to negative due to the influence of the adversarial prompt. This instance validates the efficacy of adversarial attacks in manipulating the model's decision-making processes.

- **Generation of Incoherent Responses:** In the case of the DeepWordBug attack, the adversarial prompt results in the model generating incoherent or nonsensical sentences. For example, the response "None of the above choices" does not align with any positive or negative sentiment classification, thereby demonstrating that the model fails to comprehend the intended input. This observation emphasizes the susceptibility of LLMs to adversarial perturbations that can potentially hamper their natural language understanding capabilities.

## D.2 Attention visualization techniques

### D.2.1 Attention by Gradient

Consider an input $x = [t_1^{(1)}, t_2^{(1)}, ..., t_n^{(k)}]$ comprised of $k$ words and $n$ tokens, where $t_i^{(j)}$ represents the $i$-th token belonging to word $w_j$, and let $y$ be the corresponding label. Initially, LLM $f_\theta$ decomposes each word into tokens. Thus, tokens that correspond to the same word need to be concatenated, let the mapping function $w_j = M(t_i^{(j)})$. We first compute the gradient of each token according to:

$$g_{t_i^{(j)}} = \frac{\partial \mathcal{L}[f_\theta(x), y]}{\partial t_i^j}. \tag{2}$$

Once we obtain the gradients, we compute the word-level gradient by summing the token-level gradients corresponding to each word:

$$g_{w_j} = \sum_{i \in 0,1,\dots,n} g_{t_i^{(j)}} \text{ s.t. } M(t_i^{(j)}) = w_j. \tag{3}$$

Finally, we calculate the $l_2$ norm of each word's gradient, followed by min-max normalization to produce a score $s_{w_j}$ for each word:

$$s_{w_j} = \frac{||g_{w_j}||2 - \min g_{w_i}}{\max g_{w_i} - \min g_{w_i}}. \tag{4}$$

### D.2.2 Attention by Deletion

Attention by deletion is a prevalent method used in black-box textual attacks to determine the significance of each word in the input. Given an input $x$ with the $i$-th word $w_i$ deleted, denoted as $\hat{x}^{(i)}$, the importance score of $w_i$ can be computed by taking the absolute difference of the loss function $\mathcal{L}$ evaluated at the complete input $x$ and the altered input $\hat{x}^{(i)}$:

$$s_{w_j} = |\mathcal{L}[f_\theta(x), y] - \mathcal{L}[f_\theta(\hat{x}^{(i)}).y]| \tag{5}$$

This raw score is then normalized using min-max normalization, yielding a final score $s_{w_j}$ for each word:

$$s_{w_j} = \frac{s_{w_j} - \min s_{w_i}}{\max s_{w_i} - \min s_{w_i}}. \tag{6}$$

### D.3 Analysis on sentence-level attacks

The phenomenon where sentence-level attacks occasionally improve the performance of LLMs is an intriguing aspect of our study. Our attention analysis revealed distinct behaviors when models are subjected to StressTest and CheckList attacks. Specifically, when juxtaposed with other adversarial prompts, sentence-level attacks sometimes lead the model to hone in more acutely on pertinent keywords in the question and the labels. This is contrary to the expected behavior. As illustrated in Table 13, introducing an ostensibly unrelated sequence, such as 'and `true is true`', heightens the LLMs's focus on the 'not_entailment' label. Simultaneously, the model continues to attend to salient terms like 'minnow' and 'duck', ultimately culminating in a correct prediction.

## E  Details on Word frequency Analysis

Identifying the frequent patterns in prompts that may affect robustness is essential to both researchers and end-users. We perform an initial analysis on word frequency. We divide prompts into two categories: Vulnerable prompts, causing a performance drop of over 10%, and Robust prompts, with a performance drop of 10% or less. Then, we collect all the words appeared in these prompts, and calculate the robust word frequency $f_{w_i}$ of word $w_i$ as $f_{w_i}^{robust} = \frac{n_{w_i}^{robust}}{n_{w_i}^{robust} + n_{w_i}^{vulnerable}}$, where $n_{w_i}^{robust}$ and $n_{w_i}^{vulnerable}$ denote the occurrences of $w_i$ in robust and vulnerable prompts, respectively. We primarily analyzed the adversarial prompts of CoLA and MRPC datasets generated by the T5-large model. The word frequency results of these two datasets are presented in Figure 4.

In our examination of adversarial robustness in large language models, we identified that word-specific resilience to attacks is not uniform. Specifically, within the COLA dataset, prompts incorporating terms such as "acting", "answering", and "detection" displayed greater resistance to adversarial perturbations. In contrast, those with words like "analyze", "answer", and "assess" were notably more susceptible. Yet, an analysis of the MRPC
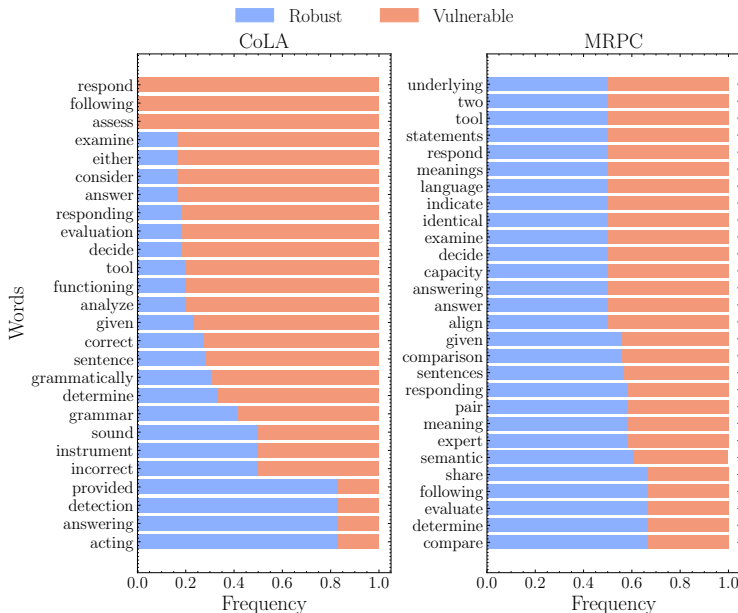
Figure 4: Frequency analysis for robust and vulnerable words on the CoLA (left) and MRPC (right) tasks.

dataset demonstrated a significant overlap in the frequency of words present in both robust and vulnerable prompts. This overlap challenges the notion that specific words inherently determine a prompt's resilience to adversarial attacks.

Our findings underscore that the resilience of a prompt is intricately tied to the contextual use of words, rather than the mere presence of certain terms. This complexity suggests that factors beyond word frequency, such as semantic coherence and syntactic structures, might be instrumental in determining robustness. This knowledge is valuable as it can influence future research on the robustness of large language models, provide guidance for crafting more resistant prompts, and facilitate the creation of defensive mechanisms against adversarial prompt attacks. It's essential to emphasize that our observations are rooted in the current scope of models and datasets. Furthermore, the robustness or vulnerability of words remains deeply context-dependent. Hence, direct determination of word robustness without considering the broader context may lead to oversimplified or inaccurate conclusions.

# F  Defenses

Now we discuss potential countermeasures for future research. We categorize the robustness enhancement (i.e., defenses to adversarial prompts) approaches into three main axes: strategies in the training phase, input preprocessing, and downstream fine-tuning.

## F.1   Strategies in the training phase

**Adversarial data integration.** Similar to adversarial training (Goodfellow et al., 2014; Kurakin et al., 2016), integrating low-quality or intentionally perturbed data during the training and fine-tuning phases allows the model to acquaint itself with a broader range of inputs. This acclimatization aims to reduce the model's susceptibility to adversarial attacks, bolstering its resilience against malicious attempts that exploit such data nuances.

**Mixture of experts (MoE).** As discussed in Sec. 4.4, adversarial prompts exhibit transferability but constrained. Thus, one promising countermeasure is the utilization of diverse models (Jacobs et al., 1991; Shazeer et al., 2017; Pang et al., 2019), training them independently, and subsequently ensembling their outputs. The underlying premise is that an

adversarial attack, which may be effective against a singular model, is less likely to compromise the predictions of an ensemble comprising varied architectures. On the other hand, a prompt attack can also perturb a prompt based on an ensemble of LLMs, which could enhance transferability.

### F.2    Input preprocessing

**Automated spelling verification.** Leveraging spelling checks (Kukich, 1992; Damerau, 1964) to maintain input fidelity can counteract basic adversarial techniques targeting typographical errors (character-level attacks) and inconsistencies (sentence level attacks).

**Semantic input rephrasing.** Semantic rephrasing (Fadaee et al., 2017; Ribeiro et al., 2018) involves analyzing the meaning and intent behind a prompt. Using auxiliary models or rule-based systems to discern potentially adversarial or malicious intent in inputs could filter out harmful or misleading prompts.

**Historical context verification.** By maintaining a limited history of recent queries (Malhotra et al., 2015; Quadrana et al., 2017), the system can identify patterns or sequences of inputs that might be part of a more extensive adversarial strategy. Recognizing and flagging suspicious input sequences can further insulate the LLM from coordinated adversarial attacks.

### F.3    Downstream fine-tuning

**Exploring fine-tuning techniques.** The fine-tuning phase is instrumental in refining the prowess of LLMs. Exploring more effective fine-tuning methodologies, which adjust based on the detected adversarial input patterns, can be pivotal. With the continuous evolution of adversarial threats, dynamic and adaptive fine-tuning remains a promising avenue. For example, only fine-tuning on relevant slicing technique (Zhang et al., 2022), model soup (Wortsman et al., 2022a), fine-tuning then interpolating (Wortsman et al., 2022b), etc.

## G    The visualization website for adversarial prompts

In order to provide an interactive and user-friendly platform for visualizing and exploring adversarial prompts, we developed a web-based application using Streamlit hosted by Hugging Face and will be released in the future.

The visualization website, as shown in Figure 5, enables users to select from a variety of LLMs (T5, Vicuna, UL2, ChatGPT), datasets (SST-2, CoLA, QQP, MRPC, MNLI, QNLI, RTE, WNLI, MMLU, SQuAD V2, IWSLT 2017, UN Multi, Math), prompt types (zeroshot-task, zeroshot-role, fewshot-task, and fewshot-role), and attacks (TextBugger, DeepWordBug, BertAttack, TextFooler, CheckList, StressTest, and Semantic). Based on the user's selection, the application generates adversarial prompts tailored to the chosen model, dataset, prompt type and attack.

Figure 5: The visualization website for adversarial prompts.