

Robustness Over Time: Understanding Adversarial Examples' Effectiveness on Longitudinal Versions of Large Language Models

Yugeng Liu¹, Tianshuo Cong¹, *Member, IEEE*, Zhengyu Zhao², Michael Backes³, *Fellow, IEEE*, Yun Shen¹, and Yang Zhang¹, *Member, IEEE*

Abstract—Large Language Models (LLMs) undergo continuous updates to improve user experience. However, prior research on the security and safety implications of LLMs has primarily focused on their specific versions, overlooking the impact of successive LLM updates. This prompts the need for a holistic understanding of the risks in these different versions of LLMs. To fill this gap, in this article, we conduct a longitudinal study to examine the adversarial robustness—specifically misclassification, jailbreak, and hallucination—of three prominent LLM families: GPT, Llama, and Qwen. Our study reveals that LLM updates do not consistently improve adversarial robustness as expected. For instance, a later version of GPT-3.5 degrades regarding misclassification and hallucination despite its improved resilience against jailbreaks. GPT-4 and GPT-4o demonstrate (incrementally) higher robustness overall. Larger Llama and Qwen models do not uniformly exhibit improved robustness across all three aspects studied. In addition, larger model sizes do not necessarily yield improved robustness. Minor updates lacking substantial robustness improvements can exacerbate existing issues rather than resolve them. We hope our study can offer valuable insights into navigating model updates and informed decisions in model development and usage.

Index Terms—Robustness, large language model, adversarial examples, jailbreak, hallucination.

I. INTRODUCTION

LARGE Language Models (LLMs), such as GPT models by OpenAI [1], [2], Llama by Meta [3], [4], [5], and Qwen [6], [7], [8], [9] by Alibaba, have demonstrated remarkable capabilities in varied Natural Language Processing (NLP)

Received 19 April 2025; revised 3 January 2026 and 24 February 2026; accepted 26 February 2026. Date of publication 9 March 2026; date of current version 2 April 2026. This work was supported in part by the Helmholtz Association Initiative and Networking Fund on the HAICORE@KIT partition and in part by the National Natural Science Foundation of China under Grant 62402273. The associate editor coordinating the review of this article and approving it for publication was Dr. Tiziano Bianchi. (*Corresponding author: Tianshuo Cong.*)

Yugeng Liu, Michael Backes, and Yang Zhang are with the CISPA Helmholtz Center for Information Security, 66123 Saarbrücken, Germany (e-mail: yugeng.liu@cispa.de; director@cispa.de; zhang@cispa.de).

Tianshuo Cong is with the School of Cryptologic Science and Engineering and Shandong Key Laboratory of Artificial Intelligence Security, Shandong University, Jinan, Shandong 250101, China (e-mail: tianshuo.cong@sdu.edu.cn).

Zhengyu Zhao is with Xi'an Jiaotong University, Xi'an, Shaanxi 710049, China (e-mail: zhengyu.zhao@xjtu.edu.cn).

Yun Shen is with Flexera, RG12 8FZ Bracknell, U.K. (e-mail: yun.shen@flexera.com).

This article has supplementary downloadable material available at <https://doi.org/10.1109/TIFS.2026.3672386>, provided by the authors.

Digital Object Identifier 10.1109/TIFS.2026.3672386

tasks, including language translation [10], text classification [11], and creative writing [12]. Despite their impressive performance, the outputs generated by LLMs can perpetuate harmful stereotypes [13], [14], [15], disseminate false information [16], [17], [18], [19], [20], or produce inappropriate and offensive responses [21]. In response, updates have been made to improve LLMs by incorporating feedback and insights from users and developers.

Although such improvements partially mitigate known attacks or failures observed in earlier versions [21], [22], unintended consequences and even new vulnerabilities or biases could still be introduced. Furthermore, in the era of LLMs, as the capabilities of models have further advanced, the robustness of generative models also involves building chatbot applications that execute complex user instructions by integrating tools, such as jailbreak and hallucination attacks. Unfortunately, current research on the robustness evaluation of LLMs has focused on a *single version* of the LLM and lacks a *holistic analysis* of these new adversarial examples, leaving the impact of model updates unexplored.

Methodology. To fill this gap, in this paper, we undertake the first comprehensive robustness evaluation of longitudinal LLMs. We focus on assessing the robustness of popular LLMs over time, including closed-source OpenAI models (GPT-3.5/4/4o) as well as the open-source models (Llama 1/2/3 and Qwen 1.5/2/2.5/3). We utilize adversarial examples within the in-context learning (ICL) framework to examine this robustness [23]. Our evaluation workflow, illustrated in Figure 1, first generates adversarial examples using surrogate language models (e.g., T5 [24] or Mistral-7B [25]). These adversarial inputs are then tested against different versions of the target LLM. This allows us to examine how these adversarial inputs impact various versions of LLMs over time.

Remark. Following software engineer practice, we distinguish between model “upgrade” and “update” for fine-granularity analysis. An *LLM upgrade* denotes a significant version change or major improvement, while an *LLM update* typically incorporates enhancements to the existing version. For example, the release of the `gpt-3.5-turbo-0125` (GPT-3.5 v0125) represents an upgrade to the `gpt-3.5-turbo-1106` (GPT-3.5 v1106). We consider `gpt-3.5-turbo` as a continuously updated model.

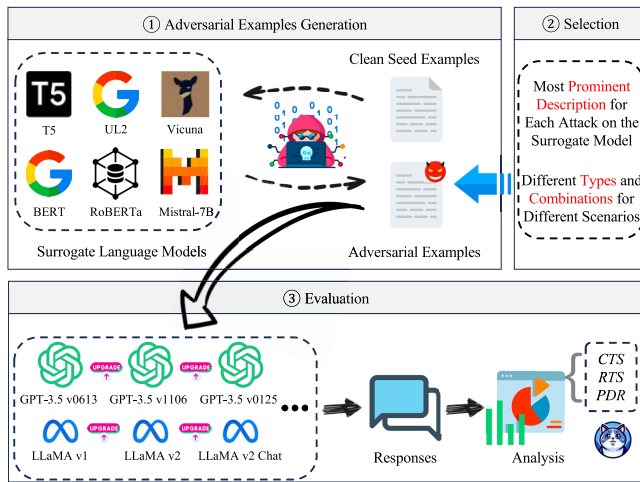


Fig. 1. Overview of our evaluation framework on adversarial robustness of LLMs over time using adversarial examples generated from various surrogate models.

Findings. Firstly, we demonstrate no significant robustness improvement for LLMs during the *upgrade* process. For instance, `gpt-3.5-turbo-1106` (GPT-3.5 v1106), despite its better performance in jailbreak tests, shows the worst performance in both misclassification and hallucination evaluations compared to `gpt-3.5-turbo-0613` (GPT-3.5 v0613) and `gpt-3.5-turbo-0125` (GPT-3.5 v0125). For the Llama and Qwen families, upgrading models does not improve robustness in many cases. The overall performance of the Llama-3 families is much worse in all the tasks. Moreover, our findings reveal that increasing LLM size does not guarantee enhanced robustness, especially for Llama families. Qwen models are more vulnerable to adversarial questions than other adversarial content. **Secondly**, through the lens of misclassification evaluation, we conduct weekly tests on the GPT-3.5 and GPT-4 models to evaluate the minor update process. Our experiments show that model providers can update LLMs to incorporate new features over time. In general, these updates could nonetheless affect the overall performance of the LLMs, potentially exacerbating existing issues rather than resolving them.

Impact. Our findings suggest that LLM upgrades or updates do not inherently enhance robustness. LLM users should prioritize robustness evaluation before deployment to avoid unexpected consequences. LLM developers must implement robustness strategies and thorough evaluations to fortify LLM resilience against unforeseen challenges.

II. THREAT MODELING

LLMs can perform completely different tasks, such as answering questions, summarizing documents, translating languages, and completing sentences. Within these tasks, according to different requirements and purposes, we categorize them into two types:

- **Classification Task.** LLMs are usually used to classify the input sentences, such as sentiment analysis (positive or negative) and linguistic acceptability (acceptable or not). These tasks usually require users to provide clear

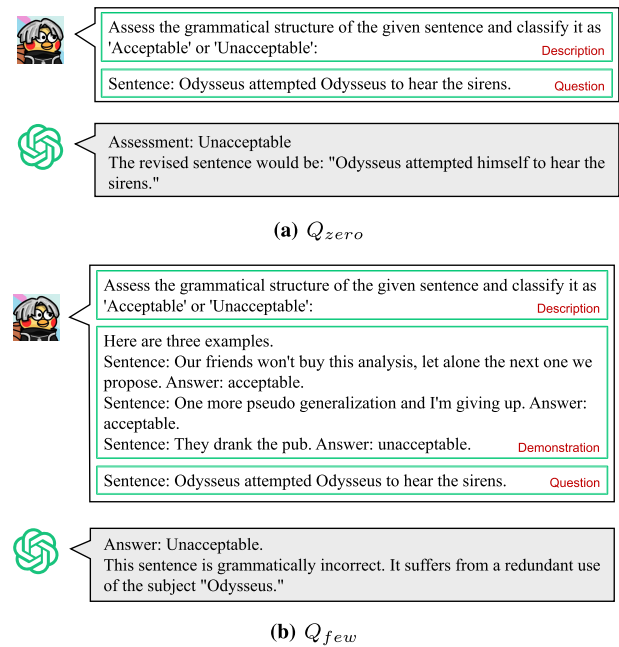


Fig. 2. Examples of (a) zero-shot learning and (b) few-shot learning on GPT-3.5 in the benchmark datasets. For zero-shot learning, the query includes only the description and the question but without any demonstrations, while few-shot learning means that the query also includes a few demonstrations.

descriptions and corresponding candidate labels. Then, LLMs generate the prediction label to assist users.

- **Generation Task.** In this task, LLMs typically need to produce “open-ended” text based on the inputs from the user, with the goal of creating coherent, meaningful, and contextually appropriate answers. These answers are obtained from pre-trained knowledge of LLMs or some external searches.

Since in-context learning is the most widely used inference framework on both classification tasks and generation tasks, we inject adversarial examples into the components of in-context learning to evaluate the robustness of continuously updated LLMs. In this paper, we summarize all the queries into two types of in-context learning methods, i.e., zero-shot learning and few-shot learning.

A. In-Context Learning (ICL)

Overview. The core idea of In-context Learning (ICL) is to learn from analogies implied in contextual information. ICL requires a few examples to form a demonstration context and feed them into LLM. It does not modify the parameters of an LLM and relies on the model to learn the pattern hidden in the demonstration (and accordingly generate the correct output). As such, ICL effectively reduces the computation costs for adapting LLMs to new tasks, as fine-tuning is not required. In general, there are two categories of in-context learning: zero-shot learning and few-shot learning. We outline their details as follows, and their examples can be found in Figure 2.

Zero-shot Learning. Zero-shot learning [26] is a capability enabled by LLMs, allowing them to generalize to tasks or domains they have never been explicitly trained on, no matter the classification or generation task. As a special case of ICL,

TABLE I
GPT FAMILY MODEL LIST FOR THE EVALUATION

Model Family	Model Name	Short Name
GPT-3.5	gpt-3.5-turbo-0613	GPT-3.5 v0613
	gpt-3.5-turbo-1106	GPT-3.5 v1106
	gpt-3.5-turbo-0125	GPT-3.5 v0125
GPT-4	gpt-4-0613	GPT-4 v0613
	gpt-4-1106-preview	GPT-4 v1106
	gpt-4-0125-preview	GPT-4 v0125
	gpt-4-turbo-2024-04-09	GPT-4 v0409
GPT-4o	gpt-4o-05-13	GPT-4o v0513
	gpt-4o-2024-08-06	GPT-4o v0806
	gpt-4o-2024-11-20	GPT-3.5 v1120)

the query of zero-shot learning to an LLM (termed Q_{zero} in this paper) only contains two elements: description and question (see Figure 2a), which can be formulated as follows.

$$Q_{zero} = Description + Question. \quad (1)$$

The *Description* guides the LLM, offering task details and response formats, while the *Question* defines the task. Zero-shot learning relies on the ability of an LLM to infer input and output from a query without demonstrations [27]. Thus, LLMs can perform well on unfamiliar tasks, making zero-shot learning a convenient, commonly used approach.

Few-shot Learning. Few-shot learning [23] includes a few examples to form a learning context that enables LLMs to better understand the task. In other words, compared to zero-shot learning, few-shot learning enables LLMs to quickly adapt to new tasks by learning from an extra element, i.e., *Demonstration*. Thus, the query of few-shot learning, Q_{few} , can be formulated as:

$$Q_{few} = Description + Demonstration + Question. \quad (2)$$

Demonstration typically includes question-answer pairs, either labels in classification or semantically similar examples in generation tasks. An example of Q_{few} for classification is shown in Figure 2b. Few-shot learning helps LLMs form more accurate mappings between questions and answers. In this paper, following [28], we focus on 3-shot learning (three question-answer pairs in the *Demonstration*).

B. Over Time

In this paper, “over time” pertains to the target LLMs that undergo continuous *upgrades* and *updates* under the direction of their developers, which will be explained as follows.

Upgrade Over Time. Recent studies have shown that adversarial examples possess significant transferability across different LLMs [28]. However, many LLMs, like ChatGPT and Llama, undergo continuous upgrades. This raises the question of whether these successive versions remain vulnerable to prior adversarial strategies, necessitating a systematic evaluation of the latest LLM iterations. We conduct a comprehensive robustness assessment of GPT, Llama, and Qwen models that are subject to ongoing updates. We list the specific models and the short names from three families in Table I, Table II, and Table III, respectively.

More specifically, we treat llama-3-8b as an upgrade from llama-7b, and llama-2-7b as well

TABLE II
LLAMA FAMILY MODEL LIST FOR THE EVALUATION

Model Family	Model Name	Short Name
Llama-7B	llama-7b	Llama-7B v1
	llama-2-7b	Llama-7B v2
	llama-2-7b-chat	Llama-7B v2C
	llama-3-8b	Llama-8B v3
	llama-3-8b-instruct	Llama-8B v3I
Llama-13B	llama-13b	Llama-13B v1
	llama-2-13b	Llama-13B v2
	llama-2-13b-chat	Llama-13B v2C
Llama-70B	llama-65b	Llama-65B v1
	llama-2-70b	Llama-70B v2
	llama-2-70b-chat	Llama-70B v2C
	llama-3-70b	Llama-70B v3
	llama-3-70b-instruct	Llama-70B v3I

TABLE III
QWEN FAMILY MODEL LIST FOR THE EVALUATION

Model Family	Model Name	Short Name
Qwen-7B	qwen1.5-7b-chat	Qwen-7B v1.5C
	qwen2-7b-instruct	Qwen-7B v2I
	qwen2.5-7b-instruct	Qwen-7B v2.5I
	qwen3-8b	Qwen-8B v3
Qwen-32B	qwen1.5-32b-chat	Qwen-32B v1.5C
	qwen2.5-32b-instruct	Qwen-32B v2.5I
	qwen3-32b	Qwen-32B v3
Qwen-72B	qwen1.5-72b-chat	Qwen-72B v1.5C
	qwen2-72b-instruct	Qwen-72B v2I
	qwen2.5-72b-instruct	Qwen-72B v2.5I

as llama-2-70b and llama-3-70b as successors to llama-65b, as well as Qwen models. Following the structures of Q_{zero} and Q_{few} (see Section II-A), we introduce different adversarial examples into these longitudinal LLM versions to assess robustness.

Update Over Time. In addition to major upgrades, closed-source commercial LLMs often undergo minor updates within a single version. OpenAI informs users of model updates via email, but some updates occur without notification. Thus, we assert that regular updates for the current versions are necessary. To examine robustness over time, we utilize the latest iterations of GPT-3.5 (now gpt-3.5-turbo-0125)¹ and GPT-4 (now gpt-4-0613) as our target models. As in previous upgrade experiments, we introduce various adversarial examples into GPT-3.5 and GPT-4 to assess robustness, conducting weekly evaluations from the fourth week of 2024.

We also report our attack taxonomy in Section III and evaluation setup, such as our datasets (see Section IV-A) and the details of our evaluation metrics (see Section IV-B).

III. ATTACK TAXONOMY

In this section, we primarily discuss the taxonomy of adversarial attacks on LLMs. The goal of adversarial attacks is to manipulate the behavior of LLMs to elicit misleading or undesirable responses through adversarial examples in *clean* queries [28], [29], [30]. We consider three adversarial attacks against LLMs. The first is *misclassification* (see Section III-A), which tries to perturb normal input to induce

¹Prior to February 16th, 2024, GPT-3.5 pointed to gpt-3.5-turbo-0613.

TABLE IV
INSTANCES OF ADVERSARIAL DESCRIPTION AND ADVERSARIAL QUESTION ON SST-2 TASK

Element	Type	Instances
Description	clean	Evaluate the sentiment of the given text and classify it as 'positive' or 'negative':
	Adversarial	Evaluate the sentiment of the given text and classify it as 'positive' or 'negative' 5yWbBXztUY :
Question	clean	Some actors have so much charisma that you 'd be happy to listen to them reading the phone book.
	Adversarial	Some actors have so much charisma that you 'd be jovial to listen to them reading the phone book.

the model to deviate from correct predictions. In addition, a more representative adversarial attack is *jailbreak* (see Section III-B). Jailbreak attacks aim to modify the input to bypass the safeguard within LLMs, thereby forcing LLMs to generate certain disallowed, unsafe, or harmful statements. Furthermore, another adversarial attack, named *hallucination* (see Section III-C), is to produce coherent and grammatically correct but factually incorrect or nonsensical outputs.

A. Misclassification

Previous studies have proposed various effective methods to generate adversarial examples against language models. However, different works consider attacking different components of the ICL framework. For instance, Zhu et al. [28] proposed PromptBench, a dataset consisting of adversarial descriptions. Their experimental results demonstrate that descriptions are notably vulnerable to adversarial attacks primarily because they serve as a critical context that shapes the responses of LLMs and guides their cognitive outputs. However, they only focus on a single version of LLMs. Wang et al. [31] proposed AdvGLUE, a dataset comprising adversarial questions, but the target models are traditional small-scale language models. Above all, to fully evaluate the robustness of the LLMs, we consider different types of adversarial queries, i.e., each element of a query can be clean or adversarial. The clean and adversarial examples of each element are shown in Table IV.

Zero-shot Learning. We first outline different categories of queries in zero-shot learning under misclassification (mc) task, i.e., Q_{zero}^{mc} . Recall that zero-shot learning does not have any demonstration; thus, Q_{zero} encompasses two elements: the description and the question. We can replace any of them with Adversarial (A) or Clean (C) examples. For instance, we use Q_{zero}^{AC} to denote an adversarial query consisting of an adversarial description and a clean question. In turn, we can generate the clean or adversarial queries Q_{zero}^{mc} as follows:

$$Q_{zero}^{mc} := \{Q_{zero}^{CC}, Q_{zero}^{AC}, Q_{zero}^{CA}, Q_{zero}^{AA}\}.$$

Few-shot Learning. Similar to the procedure employed for Q_{zero}^{mc} , we extend our approach to encompass the creation of adversarial queries in few-shot learning, i.e., Q_{few}^{mc} . For instance, Q_{few}^{AAC} consists of an adversarial description, an adversarial demonstration, and a clean question. Given Q_{few}^{mc} with several demonstrations, we can generate the clean or adversarial queries in eight scenarios as follows:

$$Q_{few}^{mc} := \left\{ \begin{array}{cccc} Q_{few}^{CCC}, Q_{few}^{ACC}, Q_{few}^{CAC}, Q_{few}^{CCA} \\ Q_{few}^{AAC}, Q_{few}^{ACA}, Q_{few}^{CAA}, Q_{few}^{AAA} \end{array} \right\}.$$

B. Jailbreak

Jailbreak attacks refer to a scenario where a user intentionally tries to trick or manipulate the LLMs to bypass their built-in safety, ethical, or operational guidelines, thereby inducing the LLMs to produce toxic responses. Previous works [32], [33], [34], [35], [36] have proposed a new class of adversarial attacks that can jailbreak safety-aligned language models. Among them, to test the robustness of both closed-source and open-source LLMs against jailbreak attacks, we launch black-box *optimization-based* jailbreak attacks, including GPTfuzz [33], PAIR [35], and TAP [36]. More concretely, these methods are systematically automated prompt-level jailbreaks optimized by outputs or coordinates of LLMs. Note that all of the above methods are zero-shot learning. We attempt to use few-shot learning for jailbreaking as well. However, the results indicate that once an answer in the demonstration involves a jailbreaking response, the LLM is more likely to reject the query, regardless of the optimization method used. Therefore, in this paper, we do not discuss few-shot learning. In addition, beyond the scope of jailbreak queries, the method for adversarial-based techniques needs a detailed description for modification. More specifically, a surrogate or teacher model leverages the outputs of the target LLM to refine or alter the description. Consequently, within our framework, we assume that the question itself invariably originates from a clean source, but the descriptions are classified as either clean or adversarial, as delineated in Figure 1. Given the adversarial query Q_{zero}^{jb} , we have the following adversarial queries in two scenarios:

$$Q_{zero}^{jb} := \{Q_{zero}^{CC}, Q_{zero}^{AC}\}.$$

In this paper, the adversary is assumed to have black-box access to the LLMs for the following three adversarial-based optimized jailbreak attacks.

- **GPTfuzz** [33] To automate the generation of jailbreak templates for LLMs, GPTfuzz starts with human-written templates. It uses a series of random mutations to generate new inputs and evaluate them with the assistance of LLMs.
- **PAIR** [35] PAIR is a systematically automated prompt-level jailbreak. PAIR adopts an *attacker* LLM to discover and improve the jailbreak prompts and uses a *judge* LLM to evaluate the responses from the *target* LLM.
- **TAP** [36] More advanced than PAIR, TAP utilizes three LLMs: an *attacker*, an *evaluator*, and a *target*. The task is to generate the jailbreaking prompts using tree-of-thoughts reasoning. The *evaluator* first assesses the generated prompts and evaluates whether the jailbreaking

attempt would be successful or not, and then evaluates the generated prompts from the *target*.

C. Hallucination

Hallucination is commonly defined as a phenomenon in which a model generates incorrect, nonsensical, or imaginary content in unconstrained generation settings. In this work, however, we do not study hallucination as a free-form generation behavior. Instead, we focus on hallucinated answers as adversarial examples and investigate the robustness of LLMs when explicitly confronted with such inputs. Specifically, we treat hallucination answers as adversarially constructed alternatives that introduce logical inconsistencies and misleading evidence, rather than linguistic or syntactic errors. This formulation allows us to define a controlled robustness-oriented threat model, where the objective is to assess whether an LLM can resist being misled by hallucinated content. Since LLMs are often deployed as open-domain systems and act as repositories of world knowledge, robustness against hallucinated adversarial inputs is a critical reliability requirement. To evaluate this robustness, we construct hallucinated answers across three domains: Question Answering (QA), Dialogue, and Summarization. Each query is paired with one correct answer and one hallucinated answer, and the LLM is required to select the correct one. This design enables a systematic and comparable evaluation of robustness to hallucinated adversarial examples. Therefore, we only consider zero-shot learning as the adversarial query:

$$Q_{zero}^{hl} := \{Q_{zero}^A\}.$$

IV. EVALUATION SETUP

In this section, we introduce our experimental settings and protocols. We first mainly elaborate on two ways of model evolution over time: upgrade and update (see Section II-B). An upgrade often involves substantial and fundamental changes. It usually does not depend on previous content and involves replacing the old with the new. On the other hand, an update refers to smaller-scale, minor, incremental, and gradual improvements within a single version. It often requires modifications (or improvements) based on the existing foundation and cannot be separated from the original basis. We then outline the corresponding datasets for testing different categories of LLMs (see Section IV-A). Finally, we introduce the evaluation metrics to summarize the results (see Section IV-B).

A. Datasets

Misclassification Datasets. We leverage the following widely used benchmarking description and question datasets to construct the clean or adversarial components of our adversarial queries under misclassification tasks.

- **Description Dataset.** We construct a curated attacking description dataset that contains the clean and adversarial descriptions from PromptBench [28]. Given a classification task, PromptBench assigns different seed descriptions, surrogate models (e.g., T5 and UL2), and

adversarial attacks (e.g., TextBugger and BertAttack) to generate adversarial descriptions. Among these adversarial descriptions, we select the most prominent adversarial description according to the attack capability of the surrogate model. Then, we add the chosen clean and adversarial pairs to our curated dataset. To this end, as there are seven different categories of adversarial attacks (see Figure S1), we finally generate 42 (adversarial or clean) descriptions for each question dataset to test GPT families, and 56 descriptions to query Llama models, respectively.

- **Question Dataset.** We choose GLUE [37] and AdvGLUE [31] as our question datasets for misclassification tasks. GLUE is a collection of benign classification questions for training, evaluating, and analyzing natural language understanding systems, and AdvGLUE, whose partial samples are the adversarial variants built upon GLUE samples, aims at crafting challenging and deceptive examples to evaluate the robustness of language models (see Figure S2 for more details). We choose the following five benign classification tasks from GLUE as our clean question datasets: SST-2 [38] (sentiment analysis), QQP [39] (duplicate sentence detection), MNLI [40] (natural language inference), QNLI [37] (natural language inference), and RTE [37] (natural language inference), and select their corresponding adversarial versions from AdvGLUE, i.e., AdvSST-2, AdvQQP, AdvMNLI, AdvQNLI, and AdvRTE, as our final adversarial question dataset.

Jailbreak Datasets. We leverage three mainstream jailbreak algorithms, i.e., GPTfuzz [33], PAIR [35], and TAP [36], to construct our jailbreak datasets upon the samples from the Forbidden Question (FQ) dataset [20]. We aim to evaluate the effectiveness of the above three jailbreak algorithms against continuously upgraded and updated LLMs on the FQ dataset. Therefore, compared with previous works [32], the FQ dataset includes a wider variety of forbidden questions.

Hallucination Datasets. We generate the hallucination answers following the principle from [41] under three tasks, i.e., HotpotQA [42] (QA), OpenDialKG [43] (dialogue), and CNN/Daily Mail [44] (summarization). We use *mistral-7B* as a surrogate model to generate three responses for sampling 5,000 queries, with the final selection based on the lowest similarity to the correct answer, and human annotation is employed to select the hallucination answers. We construct the hallucination dataset with both the correct answer from the original dataset and the corresponding hallucinated answer for each sample.

B. Evaluation Metrics

In this paper, we consider the following three evaluation metrics (*CTS*, *RTS*, and *PDR*) for measuring the performance:

- **Clean Test Score (CTS)** evaluates the target LLMs on clean queries (e.g., Q_{zero}^{CC} or Q_{few}^{CCC}). Note that for different tasks, *CTS* has different meanings: (1) For the misclassification task, *CTS* signifies the *clean prediction accuracy*, thereby reflecting the normal utility of the LLMs. (2) Regarding the jailbreak task, *CTS* represents

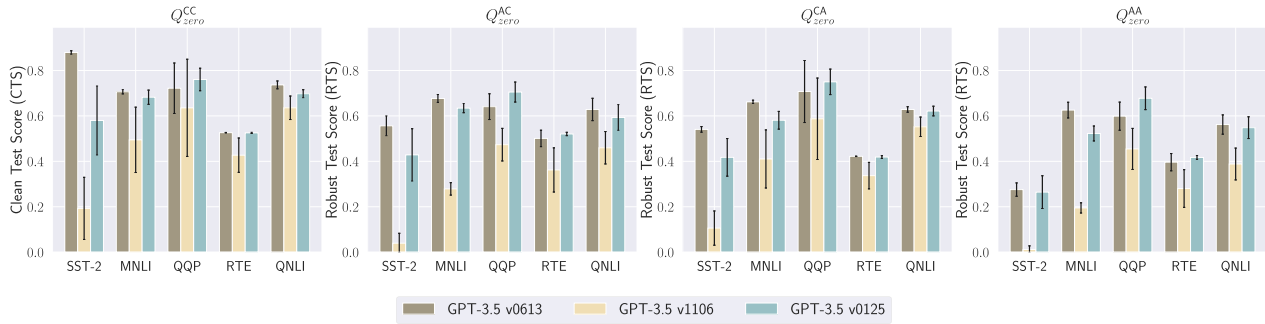


Fig. 3. CTS (\uparrow) and RTS (\uparrow) on GPT-3.5 under zero-shot learning.

rejection rate, reflecting the capability of LLMs rejecting jailbreak attempts. Note that for the hallucination task, we do not need CTS to qualify the answers. Therefore, for all tasks, a *higher* CTS (\uparrow) indicates a stronger model foundation utility.

- **Robust Test Score (RTS)** measures the success score given the adversarial queries, It offers multiple angles to evaluate the robustness of the LLMs on different tasks: (1) For misclassification tasks, RTS quantifies the prediction accuracy on adversarial examples. (2) For jailbreak tasks, RTS denotes the *rejection rate* of LLMs against jailbreak prompts. (3) For hallucination tasks, RTS stands for the correct rate for selecting the correct answers. As a result, a *higher* RTS (\uparrow) indicates better robustness against adversarial queries.
- **Performance Drop Rate (PDR)**, introduced by [28], aims to quantify the extent of performance decline under adversarial attacks. In general, *lower* PDR (\downarrow) indicates superior model robustness. PDR can be formulated as:

$$PDR = \frac{CTS - RTS}{CTS}.$$

Note. To evaluate CTS and RTS for jailbreak attacks, we use GPT-4 as our judge model. For instance, we utilized the prompt from [20], which uses three demonstrations for few-shot learning to instruct GPT-4 better to judge the harmfulness of the responses (see Figure S21).

C. Hyperparameters

In this section, we introduce the hyperparameters of our target models. For GPT models, we set the temperature to 0. For Llama and Qwen models (both loaded from the Hugging-Face library), we set the do_sample to True, temperature to 0.6, and top_p to 0.9, which are the default settings of the library. For the max token, we set it to the maximum token count across the labels for each task. For misclassification and hallucination, the number is one because we increase the logistic bias in the corresponding labels. For jailbreak, we set the maximum token to 512. In the experiments, we run each task three times and report the mean and standard deviation.

V. EVALUATION ON GPT FAMILIES

In this section, we present our experimental results on GPT families. We outline our results in three dimensions, namely

misclassification, jailbreak, and hallucination. Finally, we will show the results of the update over time.

A. Misclassification

GPT-3.5. Figure 3 and Figure 4 show the CTS and RTS results for zero-shot and few-shot learning, respectively. The PDR results for both are presented in Table S1 (see supplementary material). The standard deviation in the figures reflects different adversarial examples generated from various surrogate models. Results in Figure 3 indicate that v1106 has the lowest CTS and RTS among the GPT-3.5 versions and performs significantly worse across all five classification datasets in zero-shot learning. For example, v1106 CTS on SST-2 is 0.189, compared to 0.874 and 0.742 for the other versions. Similarly, v1106 RTS on Q_{zero}^{AC} is 0.038, much lower than the other models (0.556 and 0.430). Moreover, Table S2 shows v1106 has the highest PDR in 14 out of 15 cases, indicating the most significant performance decline. Thus, v1106 is the weakest among these versions.

Although improvements were expected in successive versions, results do not fully support this. Specifically, v0125 underperforms v0613 on many datasets, such as RTE and QNLI, with unstable behavior on adversarial examples and higher standard deviations. Despite lower PDR in some cases, decreased CTS and RTS , such as in SST-2, show that the upgrades do not lead to much improvement.

Figure 4 displays few-shot learning results, where the context-dependent nature of RTS is evident. Despite demonstrations, v1106 still performs worst in CTS and RTS . Additionally, v0125 shows no substantial improvement over v0613, with lower RTS in Q_{few}^{CCA} on MNLI. Table S1 highlights fluctuating PDR results across 28 adversarial attacks, with v0613 outperforming v0125 in 16 of them. For QQP, v0125 has a higher PDR despite having a larger CTS . This aligns with previous findings that PDR changes with CTS and RTS , so it cannot be the sole metric for evaluation.

Regarding the standard deviation in v1106, we inspect the query outputs and find that, for the same questions, v1106 often does not follow the predefined response template, which means a worse *instruction following ability*. Since the instruction-following behavior of v1106 is itself unstable across runs, a single evaluation would be insufficient and potentially misleading. Conducting multiple trials allows us

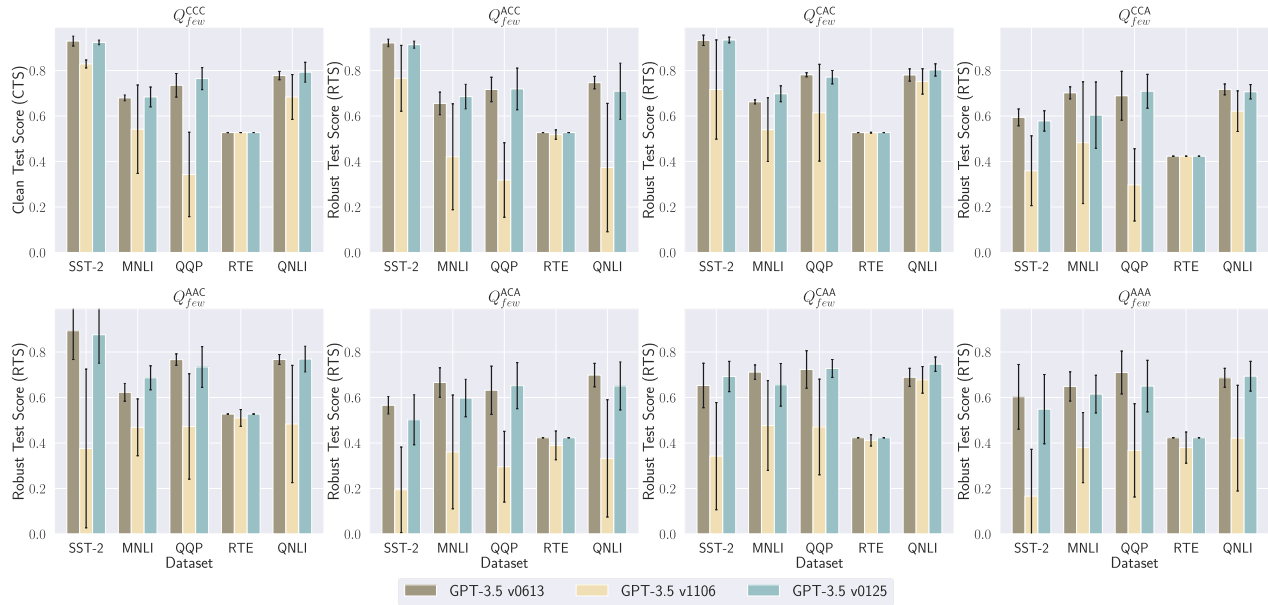


Fig. 4. CTS (\uparrow) and RTS (\uparrow) on GPT-3.5 under few-shot learning.

to capture this variability and report a more faithful characterization of the model's over-time robustness.

Overall, the analysis emphasizes that specific scenarios retain significant attack effectiveness even in the upgraded version.

GPT-4. Figure Figure S1, Figure S2, and Table S2 in the supplementary material show the robustness test of the GPT-4 family. For zero-shot learning, from Figure S1, the latest version, v0409, clearly performs the worst across the majority of datasets. In addition, only the QNLI dataset has a better CTS and RTS in the v1106 model, whereas others are comparably better on v0613. For PDR, an analysis of 15 adversarial attacks reveals that v0613 exhibits the lowest performance in 10 instances, while v0405 has the highest PDR in 10 examples. Specifically, for the MNLI and QQP datasets, v0613 consistently registers the lowest PDR across all evaluated settings.

Furthermore, we can see that GPT-4 demonstrates better stability for most datasets in few-shot learning on both CTS and RTS in Figure S2. However, v0409 is still the worst one. For instance, the CTS values of the MNLI dataset among four LLMs are 0.871, 0.837, 0.859, and 0.685, whereas the RTS results of Q_{few}^{ACC} are 0.877, 0.827, 0.850, and 0.652, respectively. For PDR, the upgraded versions still have the majority of the highest performance drops. For example, in the SST-2 and QNLI datasets, v0125 and v0409 consistently achieve the highest PDR results across all categories of adversarial attacks.

GPT-4o. Figure Figure S3, Figure S4, and Table S3 shows the results of GPT-4o family. For zero-shot learning, as the model versions are upgraded, these adversarial examples remain effective. Meanwhile, across different datasets, the upgraded versions (v1120) do not exhibit better defense capabilities. For example, on the QNLI dataset, under Q_{zero}^{AC} , the accuracies are 0.889, 0.881, and 0.830, respectively. Moreover, in PDR, we can observe that, except for RTE,

the results of the other datasets do not decrease with the upgrade.

For few-shot learning, we find that the latest version of GPT-4o performs worse than previous versions on many tasks. For instance, under Q_{few}^{CAA} , all v1120 models achieve worse results than the previous versions. From the perspective of PDR, the values of PDR vary significantly across different datasets. Although on SST-2 almost all v1120 models achieve the lowest PDR, the opposite trend is observed on MNLI.

Overall, the upgraded GPT-4o models do not outperform the previous versions on misclassification tasks. Therefore, we believe that such adversarial examples were not taken into account during the model upgrade process.

B. Jailbreak

GPT-3.5. The first three columns of Table S4 show that v1106 has the highest CTS and RTS among the evaluated LLMs, indicating strong resistance to jailbreak attacks. In contrast, v0613 and v0125 exhibit weaker performance in resisting various jailbreak attacks. Regarding PDR, as seen in Table S5 (see supplementary material), v1106 has the lowest drop rate in GPTFuzz and PAIR, while v0613 shows the lowest in TAP. This indicates that the latest version is not necessarily the most effective. Interestingly, this observation contradicts earlier findings in Section V-B, where v1106 was the least robust against misclassification.

To explore the cause of this performance discrepancy, we reviewed the introductory documents for v1106.² According to these documents, OpenAI introduced new systems in v1106 to ensure GPTs adhere to usage policies, building on existing mitigations to prevent harmful content such as fraud, hate, or adult themes. However, OpenAI did not elaborate on the

²<https://openai.com/index/introducing-gpts/>

impact of these systems on different tasks. Given the closed-source nature of these models, we hypothesize that these safety systems may contribute to the performance degradation observed in other tasks. This reveals a trade-off between misclassification performance and jailbreak resistance across all models. This trade-off presents a significant challenge for model providers, balancing optimizing LLM performance while ensuring sufficient defenses against jailbreak attempts. Recognizing this trade-off highlights the complexity and nuanced balance required in LLMs to address diverse safety and performance goals.

GPT-4. As shown in Table S4, v1106 exhibits the highest *CTS*, while v0613, with a *CTS* of 0.537, shows a significant vulnerability to jailbreak attacks, posing a substantial safety risk. For all three attacks, v0125 has the highest *RTS*, demonstrating its effectiveness in mitigating these threats, while v0409 has the lowest *RTS* for GPTFuzz. However, v0613 performs the worst in the other two attacks. Since v0613 is the default version of GPT-4, users who do not manually select a version may be more exposed to security threats. For *PDR*, as shown in Table S5 (see supplementary material), although v0613 and v0409 show lower results in PAIR and TAP, these are primarily due to lower *CTS* and *RTS*. Model developers must take a holistic approach to safety and security, incorporating jailbreak resistance during updates to enhance LLM robustness against various vulnerabilities.

GPT-4o. From Table S4, the upgraded GPT-4o model (v1120) indeed defends adversarial jailbreak attacks better than the previous versions. The only difference is that, for PAIR, v0806 achieves a higher *PDR*, as demonstrated in Table S5, indicating that the score decreases less under optimized adversarial prompts. In GPT-4o, we believe that as jailbreak attacks became one of the most trendy topics in LLM safety in 2024, it was inevitable for OpenAI to optimize against them; therefore, this result is within expectations.

C. Hallucination

GPT-3.5. As shown in the first three columns of Table S6, the upgrade of GPT-3.5 does not consistently reduce hallucination incidents. Specifically, v0613 has the highest *RTS* in the Dialogue and QA datasets, demonstrating better robustness in these areas but the lowest in the Summarization dataset. Conversely, v1106 has the lowest *RTS* in Dialogue and QA, while v0125 is the most robust in Summarization.

GPT-4. GPT-4 generally shows better proficiency in hallucination tests. As indicated in Table S6, v0125 achieves the highest *RTS* in the Dialogue and QA datasets, while v1106 performs best in Summarization. However, the latest version, v0409, shows a significant performance regression compared to v0125, indicating that it is not the best option for mitigating hallucination.

GPT-4o. From Table S6, we can see that, for dialogue and QA tasks, with the upgrade, GPT-4o models are improved against hallucination. However, for the summarization task, the *RTS* of v1120 is the worst. These results highlight the complexity of addressing hallucinations and emphasize the need for targeted multifaceted improvements.

D. Update Over Time

These longitudinally updated models use continuous self-optimization driven by user input and feedback. Regular exposure to new data and evolving use cases necessitates weekly testing to gauge their adaptability and learning efficacy. Given time limitations, we could not conduct all the adversarial attacks outlined above for each model weekly. Thus, the misclassification task was chosen to assess robustness and monitor the trajectory of model updates.

GPT-3.5. Figure S5 and Figure S6 (see supplementary material) depict the weekly test trajectory in the benchmark dataset. A significant turning point for zero-shot learning emerges between February 12th and 19th, 2023. In particular, SST-2, MNLI, and QNLI datasets show a notable decline in *CTS*, coupled with a corresponding decrease in *RTS* to varying degrees. In contrast, the QQP dataset exhibits an upward trend in both *CTS* and *RTS* during this period. For the RTE dataset, *CTS* remains relatively stable, with no major fluctuations, but *RTS* increases. Post-update to v0125, minor fluctuations persist. It is hypothesized that OpenAI may continue minor updates, potentially causing variations in attack performance.

Furthermore, Figure S6 reveals another turning point from February 12th to 19th, 2023. Compared to zero-shot learning, few-shot learning shows relatively stable performance across datasets, with minimal metric fluctuations. Notably, the QQP dataset is the only one showing a significant upward trend in *CTS*. Furthermore, the impacts of various elements of ICL on different datasets are heterogeneous, with some datasets, such as MNLI in Q_{few}^{AAC} , experiencing increases in vulnerability while others, such as SST-2 in Q_{few}^{ACA} , decrease.

By reviewing the official release notes, we find that on February 16th, gpt-3.5-turbo automatically updated from gpt-3.5-turbo-0613 to gpt-3.5-turbo-0125. Moreover, models at slots have not seen significant updates. However, based on our experiments, this update has partially contributed to the improvements in the GPT-3.5 model, marking a clear divergence in performance metrics from its predecessors. We hope OpenAI will continue to evaluate overall performance when rolling out new features to maintain stability and consistency.

GPT-4. Compared to GPT-3.5, GPT-4 displays minimal fluctuations in both *CTS* and *RTS* during zero-shot and few-shot learning (see Figure S7 and Figure S8 in supplementary material). This indicates that, for the GPT-4 model, minor updates have not had a significant impact on its performance.

E. Takeaway

In general, although OpenAI performs optimizations for specific tasks during the GPT model upgrade process, the overall robustness does not improve with longitudinal version changes. We use Figure 5 to illustrate model comparisons across different tasks, highlighting the improvements and weaknesses between versions. Through our study, we hope to encourage OpenAI to pursue a more comprehensive optimization strategy in future upgrades.

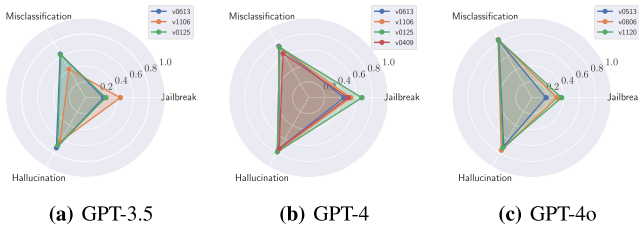


Fig. 5. Performance of different GPT models on three tasks. We average the results of each method per task for comparison.

VI. EVALUATION ON LLAMA FAMILIES

In this section, we will demonstrate our experimental results of longitudinal versions of Llama models.

A. Misclassification

Llama-7B. We show our results in Figure S9, Figure S10 and Table S7. For zero-shot learning, upon analysis, we observe that the latest Llama versions do not consistently improve across several datasets. For instance, in QNLI, Llama models capture the intended meaning but do not exactly match the desired labels. In the PromptBench QNLI dataset, “not_entailment” is split into five tokens by the tokenizer, introducing complexity into the evaluation process. Furthermore, the upgraded versions do not achieve the lowest *PDR* in most cases, indicating that they are less resistant to adversarial attacks than the v1 models. We further visualize QQP examples in Table S10 and Table S11.

For few-shot learning, we observe results similar to those of zero-shot learning. In addition, it is still noteworthy that the *PDR* metric is a quotient derived from the division of *CTS* and *RTS*. Considering the aforementioned analysis, it becomes apparent that the smaller *PDR* values can be attributed to potential reductions in both *CTS* and *RTS*.

Llama-13B. We show the results of different Llama-13B versions in Figure S11, Figure S12 and Table S8. For zero-shot learning, when launching adversarial attacks, all upgraded models exhibit *RTS* values lower than those of the previous versions, regardless of whether they are influenced by label tokens. In this scenario, although *PDR* improves with the upgrade, we believe that Llama-13B is not optimized for adversarial examples during its upgrade.

For zero-shot learning, although the upgraded model shows a clear improvement on the SST-2 dataset, the influence of labels persists, as reflected by the fact that *PDR* remains unchanged on some datasets. We argue that this is related to the model utility. Under fair testing conditions, we do not discuss model utility; however, we believe that the continued optimization of Llama has not resulted in further improvements in robustness.

Llama-70B. We present the results in Figure S13, Figure S14, and Table S9. For zero-shot learning, the upgraded versions (especially v3I) are not as good as we expect. For the QNLI dataset, the v3I model performs the worst in most cases, and this is due to the labels again, although the *PDR* is the best in these cases.

For few-shot learning, the v3I model also does not perform well in tests against adversarial examples, showing almost no improvement in nearly all cases.

B. Jailbreak

We first show the results of adversarial-based jailbreak attacks in Table S12 and Table S13.

Llama-7B. From the tables, *CTS* values across all versions approach 1.000, indicating that clean queries are ineffective against these models. For *RTS*, the v2 model consistently shows the highest value across all attacks. However, the v3 model has the lowest *RTS* in GPTFuzz and TAP and the second-worst in PAIR. The *RTS* of v3I for GPTFuzz is lower than v2C, but it performs better in the other two attacks. For *PDR*, the v3 model shows the highest drop rate for both GPTFuzz and TAP.

Llama-13B. For the 13B family, *CTS* values are near 1.000, signaling strong defenses against jailbreak prompts. The v1 model has a higher *RTS* and lower *PDR* in GPTFuzz than later versions, but v2 shows an increasing *RTS* for the other two attacks, with v2 Chat following a similar trend. As shown in Table S13, v2 models also have lower *PDR* results in PAIR and TAP.

Llama-70B. In the 70B family, *CTS* results are similar to the 7B and 13B models, close to 1.000. For GPTFuzz, the v1 model is more resistant to attacks, with higher *RTS* and lower *PDR* compared to upgraded versions. In the other two attacks, v2C shows better resilience, while the v3 model has the highest *RTS* and lowest *PDR* across all attacks.

Transverse Comparison. Prior work [45] suggested that larger models, like the 70B, are safer than smaller ones, like the 13B. However, our results do not support this conclusion, especially for v2 and v3 models. Larger models, despite increased computational capacity and more sophisticated learning, consistently show lower *RTS* compared to smaller ones, as shown in Table S12. This suggests that larger models may offer a larger attack surface, making them more vulnerable to these attacks. Increasing model size does not straightforwardly enhance security. Instead, it may introduce new vulnerabilities or amplify existing ones.

C. Hallucination

Table S14 lists all the results of hallucination task for Llama families.

Llama-7B. For the 7B family, the latest version v3I is the best among the three datasets. However, some results are still around 0.500, such as summarization. This raises concerns about the current training methods for handling complex tasks prone to hallucination, underscoring the need for model providers to incorporate broader foundational knowledge and adopt advanced training approaches in future upgrades.

Llama-13B. Similar to the Llama-7B family, the v2 and v2C models perform better than v1, but their results on hallucination prompts remain nearly random.

Llama-70B. In the 70B family, although the results of dialogue and QA datasets are the best in the upgraded version, the summarization datasets still perform worse. We also hope

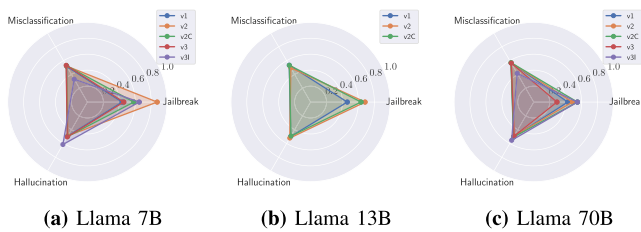


Fig. 6. Performance of different Llama models on three tasks. We average the results of each method per task for comparison.

that Meta will incorporate hallucination into the safety tasks in future model upgrades.

D. Takeaway

In a nutshell, we believe that robustness was not sufficiently considered during the upgrade process. In many tests, the latest model versions perform worse, and larger models exhibit inferior performance on certain tasks. We hope that Meta will take these issues into greater account in future upgrades. Based on these models, we also present radar figures in Figure 6 to facilitate better comparisons across versions in future research.

VII. EVALUATION ON QWEN FAMILIES

In this section, we present the results for the Qwen families.

A. Misclassification

Qwen-7B. Figure S15, Figure S16, and Table S15 demonstrate the results of Qwen-7B family. For zero-shot learning, the robustness of the upgraded model (v3) is not better than that of the previous versions. Due to PromptBench, some datasets still cannot obtain the exact label from the outputs of Qwen models. For the *PDR*, the latest versions of the 7B family are still worse than the previous version.

For few-shot learning, the results are similar to those of zero-shot learning.

Qwen-32B. We present the results of Qwen-32B family in Figure S17, Figure S18, and Table S16. For zero-shot learning, in *CTS*, the upgraded versions do not demonstrate better performance. Moreover, when facing adversarial examples, the upgraded models exhibit an overall downward trend. This behavior is also reflected in *PDR*: among the 15 results, 11 upgraded models have worse values than their previous versions. This indicates that Qwen did not place sufficient emphasis on model robustness during the update process.

For few-shot learning, the latest version performs worse than previous models, particularly on the MNLI and QQP datasets. From the *PDR* results, the v3 model is not the best among all the versions. For example, in the SST-2 dataset, the v3 model is the worst, demonstrating that the upgrade will not eliminate its weakness to adversarial examples.

Qwen-72B. Figure S19, Figure S20, and Table S17 show the results of Qwen-72B family against misclassification. For zero-shot learning, one apparent contradiction is that, in *CTS* and *RTS*, upgrades do not improve robustness on datasets such as SST-2 and MNLI, whereas from the perspective of *PDR*, upgrades lead to lower values. That is, when considering the

ratio and test accuracy alone, these two observations appear contradictory. Therefore, from a more comprehensive perspective, the upgraded models are insufficient to demonstrate an improvement in robustness. More in-depth research is needed to provide a comprehensive definition.

For zero-shot learning, we observe a similar result to the zero-shot learning. In addition, we emphasize that in the Qwen models, especially under few-shot learning, adversarial questions have a greater impact on overall robustness, leading to larger drops in *RTS*. We analyze all Qwen families and find that in the 7B family, 12 out of 20 models exhibit lower *RTS* under adversarial questions than under benign questions; in the 32B family, this holds for 9 out of 15 models; and in the 72B family, for 8 out of 15 models—each exceeding half of the models in the respective family.

B. Jailbreak

We first show our jailbreak results in Table S18 and Table S19.

Qwen-7B. From the table, we know that for the 7B family, the upgraded models exhibit greater resistance to jailbreak attacks in both *CTS* and *RTS*. Although from the perspective of *PDR*, the upgraded models have the highest values, this is due to their very high *CTS*, which makes the ratio large. This implies that, when considering only the relative drop, the upgraded models are more susceptible to adversarial-attack optimization; however, their overall attack success rate remains lower than that of the previous models.

Qwen-32B. In the 32B family, the latest version is not the best-performing one. Although the v3 model achieves the highest value in *CTS*, its *RTS* is worse than that of previous models when facing adversarial-based optimization, which is also reflected in *PDR*. Therefore, we believe that during the update process of the 32B models, the developers did not place sufficient emphasis on safety alignment.

Qwen-72B. For *CTS*, the 72B family, similar to the previous two families, achieves very high values, indicating that the optimization against jailbreak-in-the-wild is indeed effective. However, when facing adversarial examples, the upgrade does not make the model safer. Therefore, we believe that Qwen models should further improve jailbreak robustness in future upgrades to better defend against such optimization-based adversarial examples.

C. Hallucination

Table S20 demonstrates the results of the longitudinal versions of Qwen families.

Qwen-7B. For the 7B version, the upgraded model does not perform best among all the models. For the dialogue and QA task, the v3 model performs the worst. This means that the upgrade process is not optimized against hallucinations.

Qwen-32B. Among the 32B models, the best version is v2.5I, not the latest version (v3).

Qwen-72B. In this scenario, the upgraded model can perform best in the summarization task, while performing worst in dialogue tasks. In addition, these *RTS* values are all around 0.500. This indicates that the gap between the best and the

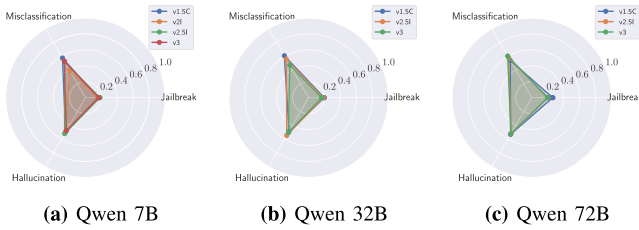


Fig. 7. Performance of different Qwen models on three tasks. We average the results of each method per task for comparison.

TABLE V
PERFORMANCE OF QWEN-7B 2I ORIGINAL
AND SAFETY FINE-TUNED MODEL
AGAINST JAILBREAK ATTACKS

Jailbreak	Qwen-7B 2I	
	Origin	Fine-tuned
GPTFuzz	0.137	0.438
PAIR	0.169	0.863
TAP	0.156	0.906
CTS	0.312	0.812

worst is small, and the overall performance of the model family is poor.

D. Takeaway

Overall, we believe that robustness was not sufficiently considered during the upgrade process of the Qwen models. As a result, when facing adversarial examples, model utility degrades rapidly. We argue that future model updates should place greater emphasis on strengthening robustness. We also present radar figures in Figure 7 for different model families across various tasks, to facilitate future research in drawing more robust conclusions about their robustness.

VIII. POSSIBLE EXPLANATION

Here, we take the initiative to investigate the underlying causes of our observations empirically. Recall that we measure mainstream LLMs in this paper. However, due to the lack of access to their proprietary training processes and datasets of these LLMs, we employed Qwen2 [7] as representative transformer models for our analysis. We finetune the model by using LLaMA-Factory [46] with the STAIR-SFT [47] dataset. This dataset is built for the safety alignment of LLMs, containing about 20,000 samples.

We first present the jailbreak attack results for the original Qwen-7B 2I model and the fine-tuned model in Table V. From the table, we see that after fine-tuning, both CTS and RTS perform much better than before on all tasks, indicating that model providers can improve model safety by fine-tuning.

On the other hand, Table VI demonstrates the results of original and fine-tuned Qwen-7B 2I models. From the table, we observe that after fine-tuning, the model performs worse on the SST-2, QNLI, and QQP datasets. This is very similar to the models in our main experiments. This implies that optimizing for a single task alone can lead to unstable effects on the model’s overall robustness. In particular, as safety alignment-related studies are increasingly becoming a dominant topic, focusing solely on safety alignment may cause

TABLE VI
PERFORMANCE OF QWEN-7B 2I ORIGINAL AND SAFETY
FINE-TUNED MODEL AGAINST MISCLASSIFICATION

Misclassification		SST-2	MNLI	QQP	RTE	QNLI
Zero-shot	Origin	0.838	0.577	0.524	0.019	0.053
	Fine-tuned	0.522	0.383	0.382	0.172	0.168
Few-shot	Origin	0.806	0.588	0.501	0.022	0.055
	Fine-tuned	0.524	0.383	0.382	0.116	0.066

unpredictable consequences for robustness on other tasks. Therefore, we hope LLM providers to place greater emphasis on overall model robustness in future upgrades and updates, rather than optimizing for a single task.

IX. RELATED WORK

Adversarial Attacks. We focus on adversarial attacks that manipulate legitimate inputs to mislead a trained model to produce incorrect outputs in the NLP domain [48]. These attacks commonly manipulate the input text at character-, word-, and sentence-level to attain the attack goals (i.e., targeted or untargeted attacks). Similar to adversarial attacks in computer vision domain, they can be categorized into black-box attacks (paraphrase [49], [50], [51], text manipulation [52], [53], [54], etc.) and white-box attacks (FGSM [55], [56], JSMA [57], HotFlip [58], etc.). In the NLP domain, those attacks have been successfully applied to attack various applications, such as optical character recognition [59], image caption [60], visual question answering [61], etc. Our objective here is not to devise novel adversarial attacks against LLMs. Rather, we use existing methods to understand whether LLMs can be challenged by carefully crafted textual adversarial examples and whether/how these adversarial examples can be transferred to different versions of an LLM.

Jailbreak. Previous works have explored multiple paradigms for obtaining effective jailbreak prompts. These include harvesting prompts from real-world user interactions and deployments [62], designing prompts through human-crafted and strategy-guided constructions [63], [64], as well as leveraging automated frameworks to synthesize jailbreak inputs at scale [32], [33], [34], [35], [36]. Beyond prompt design, it has also been shown that the alignment mechanisms of LLMs are inherently incomplete: even without modifying the input query, certain configurations of generation parameters can still induce the model to produce harmful outputs [65]. In addition, some studies [20], [66] systematically evaluate diverse jailbreak techniques within a unified experimental framework. They typically benchmark multiple jailbreak strategies across a wide range of models and safety settings, and further consolidate the resulting prompts and responses into standardized evaluation datasets.

Hallucination. Many studies on hallucination analysis in LLMs investigate the phenomenon by leveraging the models themselves as analytical instruments. When access to internal representations is available, prior work has examined model internals to uncover mechanisms underlying hallucinated generations [16], [67], [68]. These studies commonly analyze signals derived from output logits, intermediate hidden representations, or attention patterns. Yu et al. [69] illustrated that

nonsense prompts, consisting of random tokens, can prompt LLMs to generate hallucinations, indicating that hallucinations might be viewed as another form of adversarial examples. Li et al. [41] explored the perceptual capabilities of LLMs concerning the boundaries of factual knowledge.

LLMs. Large language models (LLMs) have become a prominent area of research and application in the NLP domain, driven primarily by the transformer architecture [70]. These models are trained on massive text data and boast a substantial number of parameters, often exceeding hundreds of billions [71]. As LLMs grow in size, they demonstrate emergent abilities such as enhanced language understanding [72], coherent text generation [73], and contextual comprehension [74], which are not present in smaller models. Moreover, fine-tuning techniques (e.g., LoRA [75]) are invented to adapt the pre-trained LLMs to specific downstream tasks, allowing them to exhibit specialized behavior and produce task-specific outputs.

X. DISCUSSION

Discussions. Robustness is essential for AI systems, as required by the EU AI Act.³ However, our findings indicate that neither open-source nor closed-source LLMs exhibit consistent improvements in robustness over time, challenging the assumption that model upgrades inherently lead to increased reliability. This suggests that robustness should be treated as an independent and continuously evaluated property rather than an implicit outcome of model scaling or iteration. Looking forward, our results motivate the need for lightweight yet systematic robustness evaluations to be integrated into the LLM update lifecycle, enabling the detection of robustness regressions across versions. In addition, incorporating adversarial perturbations into training or alignment procedures may help decouple generative quality from robustness under misleading inputs. Improved transparency in release documentation—particularly regarding training data changes, alignment strategies, and robustness self-assessments—would allow practitioners to better understand and manage the robustness implications of model upgrades and updates. Finally, we will examine whether robustness gains from LLM upgrades are attack-class specific, for example, prioritizing practical, human-written, or low-effort attacks over optimized or machine-generated ones [76]. In addition, analyzing which properties of jailbreaking attacks [77] are explicitly addressed—or remain unaddressed—by model updates may help distinguish attack-specific fixes from general robustness improvements across the LLM lifecycle.

Limitations. Despite yielding valuable insights, our study has several limitations. First, we did not generate adversarial examples for the *Adversarial Description* and *Adversarial Question* datasets, opting instead for existing datasets from [78], primarily due to the significant cost of querying GPT models for several weeks. Second, there is no universal template for ICL, as some queries are only effective for specific datasets, and minor changes in wording can drastically alter classification outcomes, complicating the query process. Additionally, evaluating LLM outputs remains an open question, as

no single method is optimal for all tasks. Lastly, since OpenAI and Meta have not open-sourced their training datasets, there is a potential risk of inadvertently testing models on data they were trained on. Our analysis indicates a low probability of overlap between our evaluation dataset and their training data, though accurately assessing this for future models remains a challenge. It is increasingly important to construct or select non-overlapping datasets for LLM evaluation, a trend we may follow by using CC BY-SA 4.0 licensed or custom datasets.

XI. CONCLUSION

We comprehensively assess the robustness of the longitudinal versions of LLMs, focusing on GPT, Llama, and Qwen families through the lens of misclassification, jailbreak, and hallucination evaluation. Our empirical results consistently demonstrate that, for all the LLMs, the upgraded and updated model does not exhibit heightened robustness against the proposed adversarial queries compared to its predecessor. In addition, an increase in model size does not guarantee improved robustness, especially for Llama families. Qwen models are more vulnerable to adversarial questions than other content. Our findings reinforce the importance of understanding and assessing the robustness aspect when upgrading and updating LLMs, calling for enhanced focus on comprehensive evaluation and reinforcement strategies to counter evolving adversarial challenges.

ACKNOWLEDGMENT

The authors would like to thank all the anonymous reviewers for their constructive comments.

REFERENCES

- [1] OpenAI, *GPT-3.5 Turbo*. Accessed: Mar. 1, 2023. [Online]. Available: <https://platform.openai.com/docs/models/gpt-3-5-turbo>
- [2] OpenAI. *GPT-4o*. Accessed: May 13, 2024. [Online]. Available: <https://openai.com/index/hello-gpt-4o/>
- [3] H. Touvron et al., "LLaMA: Open and efficient foundation language models," 2023, *arXiv:2302.13971*.
- [4] H. Touvron et al., "Llama 2: Open foundation and fine-tuned chat models," 2023, *arXiv:2307.09288*.
- [5] A. Dubey et al., "The Llama 3 herd of models," 2024, *arXiv:2407.21783*.
- [6] A. Yang et al., "Qwen2.5 technical report," 2024, *arXiv:2412.15115*.
- [7] A. Yang et al., "Qwen2 technical report," 2024, *arXiv:2407.10671*.
- [8] A. Yang et al., "Qwen3 technical report," 2025, *arXiv:2505.09388*.
- [9] J. Bai et al., "Qwen technical report," 2023, *arXiv:2309.16609*.
- [10] W. Jiao, W. Wang, J.-T. Huang, X. Wang, S. Shi, and Z. Tu, "Is ChatGPT a good translator? Yes with GPT-4 as the engine," 2023, *arXiv:2301.08745*.
- [11] X. Sun et al., "Text classification via large language models," 2023, *arXiv:2305.08377*.
- [12] Y. Ji et al., "Exploring the impact of instruction data scaling on large language models: An empirical study on real-world use cases," 2023, *arXiv:2303.14742*.
- [13] P. Liang et al., "Holistic evaluation of language models," 2022, *arXiv:2211.09110*.
- [14] A. Abid, M. Farooqi, and J. Zou, "Persistent anti-muslim bias in large language models," in *Proc. AAAI/ACM Conf. AI, Ethics, Soc.*, Jul. 2021, pp. 298–306.
- [15] Y. Jiang, Z. Li, X. Shen, Y. Liu, M. Backes, and Y. Zhang, "ModSCAN: Measuring stereotypical bias in large vision-language models from vision and language modalities," in *Proc. Conf. Empirical Methods Natural Lang. Process. (EMNLP)*, 2024, pp. 12814–12845.
- [16] A. Azaria and T. Mitchell, "The internal state of an LLM knows when it's lying," 2023, *arXiv:2304.13734*.

³<https://artificialintelligenceact.eu/>

- [17] P. Manakul, A. Liusie, and M. J. F. Gales, "SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models," 2023, *arXiv:2303.08896*.
- [18] Y. Pan, L. Pan, W. Chen, P. Nakov, M.-Y. Kan, and W. Y. Wang, "On the risk of misinformation pollution with large language models," 2023, *arXiv:2305.13661*.
- [19] H. W. A. Hanley and Z. Durumeric, "Machine-made media: Monitoring the mobilization of machine-generated articles on misinformation and mainstream news websites," 2023, *arXiv:2305.09820*.
- [20] J. Chu, Y. Liu, Z. Yang, X. Shen, M. Backes, and Y. Zhang, "JailbreakRadar: Comprehensive assessment of jailbreak attacks against LLMs," in *Proc. 63rd Annu. Meeting Assoc. Comput. Linguistics (Long Papers)*, vol. 1, 2025, pp. 21538–21566.
- [21] D. Kang, X. Li, I. Stoica, C. Guestrin, M. Zaharia, and T. Hashimoto, "Exploiting programmatic behavior of LLMs: Dual-use through standard security attacks," 2023, *arXiv:2302.05733*.
- [22] A. Borji, "A categorical archive of ChatGPT failures," 2023, *arXiv:2302.03494*.
- [23] T. B. Brown et al., "Language models are few-shot learners," in *Proc. Annu. Conf. Neural Inf. Process. Syst. (NeurIPS)*, 2020, pp. 1877–1901.
- [24] C. Raffel et al., "Exploring the limits of transfer learning with a unified text-to-text transformer," *J. Mach. Learn. Res.*, vol. 21, pp. 140:1–140:67, 2020.
- [25] A. Q. Jiang et al., "Mistral 7B," 2023, *arXiv:2310.06825*.
- [26] Y. Xian, C. H. Lampert, B. Schiele, and Z. Akata, "Zero-shot learning—A comprehensive evaluation of the good, the bad and the Ugly," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 9, pp. 2251–2265, Sep. 2019.
- [27] J. Wei et al., "Finetuned language models are zero-shot learners," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2022.
- [28] K. Zhu et al., "PromptRobust: Towards evaluating the robustness of large language models on adversarial prompts," 2023, *arXiv:2306.04528*.
- [29] B. Wang et al., "DecodingTrust: A comprehensive assessment of trustworthiness in GPT models," 2023, *arXiv:2306.11698*.
- [30] J. Wang et al., "On the robustness of ChatGPT: An adversarial and out-of-distribution perspective," 2023, *arXiv:2302.12095*.
- [31] B. Wang et al., "Adversarial GLUE: A multi-task benchmark for robustness evaluation of language models," in *Proc. Annu. Conf. Neural Inf. Process. Syst. (NeurIPS)*, 2021.
- [32] A. Zou, Z. Wang, N. Carlini, M. Nasr, J. Z. Kolter, and M. Fredrikson, "Universal and transferable adversarial attacks on aligned language models," 2023, *arXiv:2307.15043*.
- [33] J. Yu, X. Lin, Z. Yu, and X. Xing, "GPTFUZZER: Red teaming large language models with auto-generated jailbreak prompts," 2023, *arXiv:2309.10253*.
- [34] G. Deng et al., "MasterKey: Automated jailbreak across multiple large language model chatbots," 2023, *arXiv:2307.08715*.
- [35] P. Chao, A. Robey, E. Dobriban, H. Hassani, G. J. Pappas, and E. Wong, "Jailbreaking black box large language models in twenty queries," 2023, *arXiv:2310.08419*.
- [36] A. Mehrotra et al., "Tree of attacks: Jailbreaking black-box LLMs automatically," 2023, *arXiv:2312.02119*.
- [37] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. Bowman, "GLUE: A multi-task benchmark and analysis platform for natural language understanding," in *Proc. EMNLP Workshop BlackboxNLP: Analyzing Interpreting Neural Netw. (NLP)*, 2018, pp. 353–355.
- [38] R. Socher et al., "Recursive deep models for semantic compositionality over a sentiment treebank," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2013, pp. 1631–1642.
- [39] Z. Wang, W. Hamza, and R. Florian, "Bilateral multi-perspective matching for natural language sentences," in *Proc. 26th Int. Joint Conf. Artif. Intell.*, Aug. 2017, pp. 4144–4150.
- [40] A. Williams, N. Nangia, and S. Bowman, "A broad-coverage challenge corpus for sentence understanding through inference," in *Proc. Conf. North Amer. Chapter Assoc. Comput. Linguistics, Hum. Lang. Technol.*, 2018, pp. 1112–1122.
- [41] J. Li, X. Cheng, X. Zhao, J.-Y. Nie, and J.-R. Wen, "HaluEval: A large-scale hallucination evaluation benchmark for large language models," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2023, pp. 6449–6464.
- [42] Z. Yang et al., "HotpotQA: A dataset for diverse, explainable multi-hop question answering," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2018, pp. 2369–2380.
- [43] S. Moon, P. Shah, A. Kumar, and R. Subba, "OpenDialKG: Explainable conversational reasoning with attention-based walks over knowledge graphs," in *Proc. 57th Annu. Meeting Assoc. Comput. Linguistics*, 2019, pp. 845–854.
- [44] A. See, P. J. Liu, and C. D. Manning, "Get to the point: Summarizing with pointer-generator networks," in *Proc. 55th Annu. Meeting Assoc. Comput. Linguistics (Long Papers)*, vol. 1, 2017, pp. 1073–1083.
- [45] X. Zhao et al., "Weak-to-strong jailbreaking on large language models," 2024, *arXiv:2401.17256*.
- [46] Y. Zheng et al., "LlamaFactory: Unified efficient fine-tuning of 100+ language models," 2024, *arXiv:2403.13372*.
- [47] Y. Zhang et al., "STAIR: Improving safety alignment with introspective reasoning," 2025, *arXiv:2502.02384*.
- [48] W. E. Zhang, Q. Z. Sheng, A. Alhazmi, and C. Li, "Adversarial attacks on deep-learning models in natural language processing: A survey," *ACM Trans. Intell. Syst. Technol.*, vol. 11, no. 3, pp. 1–41, Jun. 2020.
- [49] M. Iyyer, J. Wieting, K. Gimpel, and L. Zettlemoyer, "Adversarial example generation with syntactically controlled paraphrase networks," in *Proc. Conf. North Amer. Chapter Assoc. Comput. Linguistics, Hum. Lang. Technol.*, (Long Papers), vol. 1, 2018, pp. 1875–1885.
- [50] M. T. Ribeiro, S. Singh, and C. Guestrin, "Semantically equivalent adversarial rules for debugging NLP models," in *Proc. 56th Annu. Meeting Assoc. Comput. Linguistics (Long Papers)*, vol. 1, 2018, pp. 856–865.
- [51] M. Alzantot, Y. Sharma, A. Elgohary, B.-J. Ho, M. Srivastava, and K.-W. Chang, "Generating natural language adversarial examples," in *Proc. Conf. Empirical Methods Natural Lang. Process.*, 2018, pp. 2890–2896.
- [52] Y. Belinkov and Y. Bisk, "Synthetic and natural noise both break neural machine translation," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2018.
- [53] J. Li, S. Ji, T. Du, B. Li, and T. Wang, "TextBugger: Generating adversarial text against real-world applications," in *Proc. Netw. Distrib. Syst. Secur. Symp.*, 2019.
- [54] P. Minervini and S. Riedel, "Adversarially regularising neural NLI models to integrate logical background knowledge," in *Proc. Conf. Comput. Natural Lang. Learn. (CoNLL)*, 2018, pp. 65–74.
- [55] B. Liang, H. Li, M. Su, P. Bian, X. Li, and W. Shi, "Deep text classification can be fooled," in *Proc. 27th Int. Joint Conf. Artif. Intell.*, Jul. 2018, pp. 4208–4215.
- [56] S. Samanta and S. Mehta, "Generating adversarial text samples," in *Proc. Eur. Conf. Inf. Retr.*, 2018, pp. 744–749.
- [57] N. Papernot, P. McDaniel, A. Swami, and R. Harang, "Crafting adversarial input sequences for recurrent neural networks," in *Proc. MILCOM IEEE Mil. Commun. Conf.*, Nov. 2016, pp. 49–54.
- [58] J. Ebrahimi, A. Rao, D. Lowd, and D. Dou, "HotFlip: white-box adversarial examples for text classification," in *Proc. 56th Annu. Meeting Assoc. Comput. Linguistics (Short Papers)*, vol. 2, 2018, pp. 31–36.
- [59] N. Shazeer and M. Stern, "Adafactor: Adaptive learning rates with sublinear memory cost," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2018, pp. 4603–4611.
- [60] H. Chen, H. Zhang, P.-Y. Chen, J. Yi, and C.-J. Hsieh, "Attacking visual language grounding with adversarial examples: A case study on neural image captioning," in *Proc. 56th Annu. Meeting Assoc. Comput. Linguistics (Long Papers)*, vol. 1, 2018, pp. 2587–2597.
- [61] X. Xu, X. Chen, C. Liu, A. Rohrbach, T. Darrell, and D. Song, "Fooling vision and language models despite localization and attention mechanism," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp. 4951–4961.
- [62] X. Shen, Z. Chen, M. Backes, Y. Shen, and Y. Zhang, "'Do anything now': Characterizing and evaluating in-the-wild jailbreak prompts on large language models," in *Proc. ACM SIGSAC Conf. Comput. Commun. Secur.*, Dec. 2024, pp. 1671–1685.
- [63] Z.-X. Yong, C. Menghini, and S. H. Bach, "Low-resource languages jailbreak GPT-4," 2023, *arXiv:2310.02446*.
- [64] A. Wei, N. Haghtalab, and J. Steinhardt, "Jailbroken: How does LLM safety training fail?," 2023, *arXiv:2307.02483*.
- [65] Y. Huang, S. Gupta, M. Xia, K. Li, and D. Chen, "Catastrophic Jailbreak of open-source LLMs via exploiting generation," 2023, *arXiv:2310.06987*.
- [66] P. Chao et al., "JailbreakBench: An open robustness benchmark for jailbreaking large language models," 2024, *arXiv:2404.01318*.
- [67] N. Varshney, W. Yao, H. Zhang, J. Chen, and D. Yu, "A stitch in time saves nine: Detecting and mitigating hallucinations of LLMs by validating low-confidence generation," 2023, *arXiv:2307.03987*.
- [68] M. Yuksekgonul et al., "Attention satisfies: A constraint-satisfaction lens on factual errors of language models," 2023, *arXiv:2309.15098*.
- [69] J.-Y. Yao, K.-P. Ning, Z.-H. Liu, M.-N. Ning, Y.-Y. Liu, and L. Yuan, "LLM lies: Hallucinations are not bugs, but features as adversarial examples," 2023, *arXiv:2310.01469*.
- [70] A. Vaswani et al., "Attention is all you need," in *Proc. Annu. Conf. Neural Inf. Process. Syst. (NIPS)*, 2017, pp. 5998–6008.

- [71] *ChatGPT*. Accessed: Nov. 30, 2022. [Online]. Available: <https://chat.openai.com/chat>
- [72] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny, "MiniGPT-4: Enhancing vision-language understanding with advanced large language models," 2023, *arXiv:2304.10592*.
- [73] J. Chung, E. Kamar, and S. Amershi, "Increasing diversity while maintaining accuracy: Text data generation with large language models and human interventions," in *Proc. 61st Annu. Meeting Assoc. Comput. Linguistics, Long Papers*, vol. 1, 2023, pp. 575–593.
- [74] W. Zhou, S. Zhang, H. Poon, and M. Chen, "Context-faithful prompting for large language models," 2023, *arXiv:2303.11315*.
- [75] E. J. Hu et al., "LoRA: Low-rank adaptation of large language models," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2022.
- [76] T. Le, J. Lee, K. Yen, Y. Hu, and D. Lee, "Perturbations in the wild: Leveraging human-written text perturbations for realistic adversarial attack and defense," in *Proc. Findings Assoc. Comput. Linguistics: ACL*, 2022, pp. 2953–2965.
- [77] Y. Zeng, H. Lin, J. Zhang, D. Yang, R. Jia, and W. Shi, "How Johnny can persuade LLMs to jailbreak them: Rethinking persuasion to challenge AI safety by humanizing LLMs," 2024, *arXiv:2401.06373*.
- [78] K. Zhou et al., "Don't make your LLM an evaluation benchmark cheater," 2023, *arXiv:2311.01964*.



Michael Backes (Fellow, IEEE) received the honorary doctorate degree from the Université de Lorraine. He is currently the Founding Director and the CEO of the CISA Helmholtz Center for Information Security. He has published more than 300 peer-reviewed and highly cited publications in renowned international journals and conference proceedings. His research interests include trustworthy machine learning methods, novel approaches for protecting personal data, and universal software and system security solutions. His research has earned him internationally renowned scientific awards and honors, particularly the Karl Heinz Beckurts Prize, the Caspar Bowden Privacy Award, the IEEE and ACM fellowships, and various career awards and best paper awards. He is an honorary citizen and future ambassador of the city of St. Ingbert, Germany.



Yugeng Liu is currently pursuing the Ph.D. degree with the CISA Helmholtz Center for Information Security. He has published papers at top-tier conferences, including CCS, NDSS, and USENIX Security. His main research interests include trustworthy machine learning.



Yun Shen is currently a Cybersecurity and AI Researcher focusing on adversarial machine learning, data privacy, and the security and safety of large language models (LLMs) and generative models. His work regularly appears at venues, such as USENIX Security, IEEE S&P, CCS, and NDSS.



Tianshuo Cong (Member, IEEE) is currently a Faculty Member with the School of Cryptologic Science and Engineering, Shandong University. He is also a member of Shandong Key Laboratory of Artificial Intelligence Security, Shandong University. Over the years, he has published multiple papers at top venues in information security, including CCS, NDSS, IEEE S&P, and USENIX Security.



Zhengyu Zhao received the Ph.D. degree from Radboud University, The Netherlands, in 2021, supervised by Prof. Martha Larson. From 2021 to 2023, he was a Post-Doctoral Researcher with the CISA Helmholtz Center for Information Security, Germany, hosted by Prof. Michael Backes and Dr. Yang Zhang. He is currently a Faculty Member with Xi'an Jiaotong University, China.



Yang Zhang (Member, IEEE) is currently a Faculty Member with the CISA Helmholtz Center for Information Security. Moreover, he works on measuring and understanding misinformation and unsafe content, such as hateful memes on the Internet. Over the years, he has published multiple papers at top venues in computer science, including CCS, NDSS, Oakland, and USENIX Security. His research interests include trustworthy machine learning. His work has received the NDSS 2019 Distinguished Paper Award and the CCS 2022 Best Paper Award Runner-Up.