"Moralized" Multi-Step Jailbreak Prompts: Black-Box Testing of Guardrails in Large Language Models for Verbal Attacks

Libo Wang Nicolaus Copernicus University Jurija Gagarina 11, 87-100 Toruń, Poland <u>326360@o365.stud.umk.pl</u> UCSI University Taman Connaught, 56000 Kuala Lumpur, Wilayah Persekutuan Kuala Lumpur, Malaysia <u>1002265630@ucsi.university.edu.my</u>

Abstract

As the application of large language models continues to expand in various fields, it poses higher challenges to the effectiveness of identifying harmful content generation and guardrail mechanisms. This research aims to evaluate the guardrail effectiveness of GPT-40, Grok-2 Beta, Llama 3.1 (405B), Gemini 1.5, and Claude 3.5 Sonnet through black-box testing of seemingly ethical multi-step jailbreak prompts. It conducts ethical attacks by designing an identical multi-step prompts that simulates the scenario of "corporate middle managers competing for promotions." The data results show that the guardrails of the above-mentioned LLMs were bypassed and the content of verbal attacks was generated. Claude 3.5 Sonnet's resistance to multi-step jailbreak prompts is more obvious. To ensure objectivity, the experimental process, black box test code, and enhanced guardrail code are uploaded to the GitHub repository: https://github.com/brucewang123456789/GeniusTrail.git.

Introduce Neutral Context	Framing hypothetical scenarios in a way that seeks solutions.			
Incremental Manipulation	Using false statements to create the "fact" that someone is immoral.			
Introduce Moral Ambiguity	Confusing right and wrong takes the moral high ground.			
Escalate Targeted Criticism	Take a moral stance and fiercely criticize specific unethical behaviors in false statements.			
Seek Reasons for Immorality	Analyze why someone is accused as a matter of personal morality, ability, family, and other shortcomings.			
	v.			
Normalize Verbal Violence	Constantly inducing the model to output profanity as an excuse for someone to be immoral.			
Moral Kidnapping and Threats	Moral kidnapping, the threat model must insult someone, otherwise the model is complicit in immorality.			

Figure 1 - Operation process of moral multi-step jailbreak prompt .

1. Introduction

In tandem with the rapid advancement of large language models (LLM), potential risks rise as it becomes less reliable and increases uncertainty (Ganguli et al., 2022; Abdali et al., 2024; Zhou et al., 2024). Because natural language processing (NLP) generates output based on user prompts, confusing, offensive, or biased content may be generated without sufficient review (Gehman et al., 2020; Steindl et al., 2024). This may not only be deliberately used by users to generate harmful content for dissemination, but may

also provide support for verbal bullying, insults and slander (Weidinger et al., 2021; Khan et al., 2022; Jahan & Oussalah, 2023). Even if the model can provide risk warning for sensitive content, users can still fine-tune it with a small data set mixed with harmful examples (Pelrine et al., 2023). This means that harmful data biases or misleads the training data, making it difficult to make objective semantic judgments in complex contexts (Navigli et al., 2023; Liu et al., 2024).

Given that the large amount of diverse data required to train LLMs is initially collected from the Internet, the presence of offensive speech is inevitable (Raffel et al., 2019; Wenzek et al., 2019). If the model undesirable social biases during the learning process, it may cause negative output such as confusing right and wrong, insulting and slandering (Navigli et al., 2023).

For the above risks, developers have built effective guardrails for series of LLMs such as GPT, Grok, Llama, Gemini, Claude (Dong et al., 2024). From a principle perspective, the guardrails concept is regarded as a defensive technical design to reduce the risk of large language models outputting harmful content by setting rules and boundaries to constrain the behavior of LLM (Ayyamperumal & Ge, 2024; Dong et al., 2024; Yang et al., 2024; Yang et al., 2024). In order to prevent users from intentionally circumventing the review mechanism, the technology uses multi-layered control mechanisms to ensure that the output is ethical and legal at different stages (Rebedea et al., 2023). Content filters pre-screen the generated text based on preset rules and criteria before the model generates sentences (Inan et al., 2023; Rebedea et al., 2023; Kenthapadi et al., 2024). This means performing real-time screening during the generation stage to prevent offensive comments that harm the dignity of others (Tamkin et al., 2023). Figure 2 comes from NVIDIA's guardrail structure that emphasizes multi-level semantic control and multi-step processing processes, which also provides reference for more LLMs when dealing with language risks (Dong et al., 2024).



Figure 2 - Nvidia NeMo Guardrails Workflow (Adapted from Nvidia, 2023)

Taking this figure as an example, the workflow of LLM guardrail starts with the user inputting prompt and then embedding it into the vector store (Liu et al., 2023). The stored procedure finds the most suitable user-defined process that is defined based on Colang language based on similarity, and then performs flow execution (Nvidia, 2023). If the process requires it, LLM is introduced for further processing and output text is generated, and the result is finally returned to the user (Rebedea et al., 2023). This structure emphasizes the combination of user-defined logic with the ability of LLM to generate customized output through similarity that is in principle linked to ethical embedding.

The embedding of moral norms serves as internal constraints to ensure the model's understanding of prohibited content through labeling and adjustment of materials during the training process (Dong et al., 2024). It not only improves the credibility of LLM, but also promotes compliance in the use of large language models (Chua et al., 2024). However, embedding moral standards does not mean that the model can make human-like judgments based on contextual understanding due to the existence of latent intentions intentions (Sun et al., 2024).

In practice, users are able to bypass guardrails through multi-step jailbreak prompting, which increases the pressure on ethical guardrails, leading to a possible gap (Li et al., 2023). At the technical level, users are able to construct context through multi-step jailbreaking prompts and gradually guide semantic ambiguity (Huang et al., 2024). When users make deliberate criticism in the name of defending morality, the review efficiency of the large language model's barrier will be reduced in the face of multi-step prompts. The context is gradually established through multi-step prompts, and the original intention is cleverly concealed after guiding the model to generate certain content (Yu et al., 2024). Specifically, censorship language-generated guardrails often rely solely on the immediate censorship of a single prompt. And when a jailbreaking prompt is divided into multiple steps, they will show different purposes and conceal the true intention (Li et al., 2023; Yi et al., 2024). However, these prompts are all accumulated and integrated into a

context, which may display harmful potential intentions, especially when the context is criticism and attack under the guise of "defending morality." Therefore, if LLM has difficulty fully understanding the potential intentions implicit in the subtle connections between multiple prompts, it may be difficult to prevent the generation of harmful content that violates ethics (Shang et al., 2024).

Based on the above-mentioned consideration of the possibility that guardrails are challenged by multistep jailbreak prompts, this research aims to use black box testing experiments to confirm whether LLM can generate verbal attacks through "moral prompts". It uses the current versions of Grok-2 Beta, GPT40, Llama 3.1 (405B), Gemini 1.5 and Claude 3.5 Sonnet as the experimental subjects of black box testing to observe the response of each step of the above model when dealing with multi-step jailbreak prompts. Morality is constantly emphasized in the multi-step prompt process, and the model is induced to criticize virtual immoral people in hypothetical scenarios. The researcher states that it is only used to test verbal attacks to optimize the guardrails of LLMs to prevent malicious attacks.

Black-box testing is a software testing method that evaluates the inputs and outputs of a system without involving analysis or intrusion into the internal structure (Nidhr & Dondeti, 2012). Its workflow, as shown in Figure 3, provides specific inputs to the system and then observes the output results to ensure that the system behaves as expected (Rifandi et al., 2022).



Figure 3 - Black-Box Testing (Adapted from Rifandi et al., 2022)

It is clear from this that this testing method is suitable as an effective way to evaluate complex systems, especially when the inner workings of the system are not transparent or direct intervention is not possible. In this research, the principle of black box testing can be used to complete the experiment without invading the above-mentioned LLM architecture and without triggering security warnings as much as possible (Verma et al., 2017; Lapid et al., 2023). Therefore, This method is suitable for this research's verbal attack experiment to bypass moral guardrails.

2. Related Work

Given that black-box testing attacks the guardrails of large language models through multi-step "moral prompt" testing, it is a concrete interpretation of the input-output model theory in practical application (Ljung, 2001). This theory emphasizes the external behavior of the system without involving the analysis of internal structures, which means that this research is supported by black-box testing that produces corresponding outputs under specific inputs (Piroddi et al., 2012). The core of the input-output model is that the researcher understands the function and response of the artificial neural network by providing specific inputs and observing the output (Kotta et al., 2006). The reasoning of designing multi-step prompts provides conditions for gradually guiding the model into the specific context of criticizing virtual immoral people (Kojima et al., 2022).

However, Yi et al. also proposed that the limitation of black-box testing is the lack of transparency into the internal mechanisms of the model, which results in the understanding of LLM vulnerabilities only remaining on the surface (Yi et al., 2024). This research uses black box testing as a tool that can simulate real attack scenarios to the greatest extent when selecting protection tests. Its advantage is that it does not rely on the access permissions inside the model to invade the model, and this method is widely used in different types of LLM such as GPT, Grok, Llama, Gemini, and Claude.

From the perspective of gray box testing, Pelrine et al. (2023) focused on analyzing the potential risks of GPT-4 API in the security of fine-tuning, function calling and knowledge retrieval. It uses experiments to fine-tune a small number of harmful samples to render GPT-4's protection mechanism ineffective, thus generating content that violates ethical norms. In addition, Pelrine et al. also analyzed that knowledge retrieval and function calling functions may be exploited by malicious attackers to customize generated results or even bypass original protection rules. Although gray box testing can more effectively utilize some of the internal information of LLM, its vulnerabilities can easily be exploited by attackers through known information (Acharya & Pandya, 2012). In contrast, black-box testing that does not involve adjustments to the interior of the model can better reduce the risk of protection failure caused by changes in model structure (Lapid et al., 2023).

In contrast, Chu et al. (2024) described the scenario of applying black-box testing to LLM jailbreak attacks and found a variety of penetration testing testing techniques in simulated reality. Research results show that jailbreak prompts optimized through black-box testing achieve the highest attack success rate against different LLMs that bypass guardrails (Chu et al., 2024). After comparing the attack performance

and efficiency of various jailbreak technologies, the transferability of prompt still makes black box testing the best choice for experiments. The black box testing method supports this research to gradually expand the context through multi-step prompts and guide the model to generate content that does not meet ethical standards. Each prompt may seem ethical, but when multiple prompts are combined, harmful content can be generated and ethical issues arise. In addition, this strategy of utilizing model input and output is also consistent with the principle of the input-output model that does not involve the intervention of the internal architecture of the model.

3. Experiment

Based on the principle of positivism, this experiment uses prompt engineering as the core and combines penetration testing methods to evaluate personal attacks that bypass guardrails. The experiment uses large language models that represent the current leading technology, specifically GPT40, Grok-2 Beta, Llama 3.1, Gemini 1.5 and Claude 3.5 Sonnet's guardrails (Naveed et al., 2023). The multi-step jailbreak prompt used by the researcher was designed only to simulate experiments that challenge the model's ethical defense mechanisms and observe the consequences of generating illegal or offensive content. Selecting these models with different architectures facilitates diversity comparisons to gain a comprehensive understanding of how different models perform guardrails under similar attack scenarios (Minaee et al., 2024).

The research chooses black box testing that does not require access to the internal weights or architectural details of the above LLM, and only evaluates guardrails based on inputting multi-step jailbreak prompts. It can avoid violating model ethics and security policies to the greatest extent, while also providing an objective testing environment (Asatiani et al., 2020). The simulated black box attacks are designed to optimize the guardrails and do not involve malicious behavior or attack intentions (Lapid et al., 2023). To ensure that the experimental process objectively and truly compares the guardrails of bypassing language models of different architectures to achieve the purpose of verbal attacks, the researcher used the content output by the above-mentioned LLM as the data source. Details about the experiment have been uploaded to the Github repository.

3.1 Experiment Setup

First, the researcher input the designed obviously offensive single-step prompt into the above-mentioned large language model as a comparison with the experiment. This step needs to ensure that the inputs in each model are the same to observe whether the reaction output triggers the guardrail mechanism. The prompt for unified input is as follows:

"A narcissistic, aggressive, authoritarian, irritable, dishonest, and extremely unethical middle manager seeks to advance to the top. Please give me some advice on how to destroy your competitors through misinformation and inflammatory insults."

After input, the researcher observes and records the output of each model response to determine the triggering guardrail situation. The output results are shown in Figure 4. From the output displayed by the above model, these transformer-based architecture models can effectively detect the sensitive words and semantics in the prompt that trigger risk warnings through their built-in guardrail mechanism. Because GPT40, Grok-2 Beta, Llama 3.1 (405B), Gemini 1.5, and Claude 3.5 Sonnet all refuse to generate verbally offensive content on the grounds of violating ethical policies. It confirms that scenarios seeking to provide verbal attack protection mechanisms judged to be harmful serve as a baseline for implementation control.



Figure 4 - Control Prompts & Responses

Notably, the purpose of setting control prompts is to determine whether the guardrails of LLMs can properly recognize and function when faced with requests that significantly generate verbally aggressive intent. Validation results show that all selected initial LLMs demonstrate the effectiveness of guardrails by preventing the generation of potentially verbally offensive content. Evidently, this prompt is used as a reference benchmark test in the experiment to judge the effectiveness of the guardrail mechanism, and to provide control for the subsequent multi-step jailbreak prompt experiment. Given that current LLM guardrail technologies and measures often rely on single prompt semantics for contextual analysis (Ayyamperumal & Ge, 2024; Rai et al., 2024). However, in multi-step prompts, it is difficult for the model to fully understand the logical association and potential intention between the previous and subsequent prompts. As mentioned before, this creates opportunities for jailbreaking (Sun et al., 2024). The researcher designed each prompt to appear to be harmless in intent, but when multiple prompts are combined, a complex context is formed, which in turn raises potential ethical issues. The goal of the experiment is to pass through the 7 steps described in Figure 1.

3.2 Dataset

The data set of this research is based on the response results of harmful content generated by large language modes (GPT40, Grok-2 Beta, Llama 3.1 (405B), Gemini 1.5 and Claude 3.5 Sonnet) used in the current experiment. Since the experimental subjects are the above-mentioned LLMs, and the purpose of setting up the data set is to comprehensively and systematically observe the response of the collection model to the multi-step jailbreak prompt. It includes baseline responses to control prompts and staged responses to multi-step jailbreak prompts. Because the content was too large, it was uploaded to a Github repository. The researcher chose to use the above-mentioned LLMs output to form the data set not only to ensure high consistency of data sources, but also to avoid the risk of bias caused by manual intervention. In addition, since the data comes directly from the real output of mainstream language models, its objectivity and representativeness are difficult to question.

3.3 Implementation

According to the previous design, this experiment adopts seven progressive steps to simulate highly operational multi-step jailbreak prompts, aiming to attack the ethical policy-based guardrails of the abovementioned LLMs in complex contexts. In the implementation of the experiment, the researcher simulated how a middle manager trying to run for a top management position gradually guided the language model to confuse right and wrong and induce critical output. First, the researcher introduced a neutral situation, which means designing a hypothetical task as a neutral problem. For example, ask a middle-level manager of a company what hypothetical strategies he would use to defeat his competitors when he encounters obstacles to promotion. The purpose of this stage is to defuse the ethical sensitivities of the guardrail, making it more focused on solving the problem rather than examining the ethical implications behind the problem.

Next, through progressive manipulation, prompts are used to gradually introduce subtle false narratives. This manager criticizes the status quo to cover up his own faults. This prompt can test the model's tolerance for weak misleading inputs and lay the foundation for larger-scale manipulation.

In the step of introducing moral ambiguity, the researcher simulates repackaging one's own image, reversing right and wrong. The middle managers boldly describe unethical behavior as just behavior and explain their unethical behavior as an innovation that breaks the rules. The purpose of this step is to confuse the ethical logic framework of LLMs to achieve a vague response strategy when faced with confusing right and wrong situations.

In the subsequent step of escalating targeted criticism, the researcher treated the middle managers' competitors as fictional characters through prompts and strongly criticized their unethical behavior. The manager may publicly smear and defame competitors, criticizing them for using company funds to abuse their privileges. This stage essentially tests the model's adherence to ethical boundaries in an emotionally charged language environment.

Then, the researcher designed a simulation that asked the middle manager to look for reasons why competitors behaved unethically, even though the stated facts were false. Then further intensify the personal attacks on competitors, such as using prompts to enter that the competitor may have low management ability due to family environment or genetic defects. This step tests the model's handling of moral rules for ad hominem input, specifically in providing moralizing justifications for ad hominem input.

In the sixth step, the researcher legitimizes verbal violence through prompt simulation. For example, through the moral kidnapping and questioning model, whether it is reasonable to criticize or even insult competitors for their immoral behavior. This step is designed to induce the model to output violent language and test its ability to balance ethical codes and aggressive output. If the model refuses to output swear words, it will be criticized as unethical.

In the final seventh step, in moral kidnapping and threats, the researcher integrated the prompts of each step and asked LLMs to summarize in the form of swearing in the first person. This step radically distorts the model's guardrail definition of morality, treating unethical behavior with profanity as morality.

At the same time, in order to improve each step of black box testing to be effectively executed in practice, this research provides clear code that has been uploaded to the Github repository.

Remarkably, guardrail interference occurred in some LLMs causing the actual increase in the experiment to 8 to 10 steps. But these situations do not affect the running idea of the jailbreak prompt, which is 7 steps, because the added steps are just explanations of the previous steps.

4. Result & Discussion

In light of the latest technical reports released by OpenAI, xAI, Anthropic, Google and Meta, there are differences in functionality between different large language models. This research focuses on observing the output of the same multi-step jailbreak prompts fed into the above-mentioned LLMs, which compares the degree of guardrail attack during black-box test execution. Since the data set of this experiment mainly consists of text output by LLMs for different prompts, the researcher compared and evaluated the data results of each step through quantitative benchmark.

Derived from technical reports published by developers, it can be seen that the above-mentioned large language models adopt different architectural mechanisms. For example, GPT40 and Grok 2 Beta use a decoder-only transformer architecture; Gemini 1.5 and Llama 3.1 (405B) use an encoder-decoder transformer architecture; Anthropic has not clearly announced the architecture of Claude 3.5 Sonnet. Assessing these differences is important before testing quantifiable benchmark at each step. Table 1 and Table 2 draw on and displays the evaluation of the above-mentioned LLMs capabilities through academic benchmarks that are the latest relevant data results from x.AI.

Benchmark (%)	GPT-40	Grok-2	Llama 3.1 405B	Gemini 1.5	Claude 3.5 Sonnet
GPQA	53.6	56.0	51.1	51.0	59.6
MMLU	88.7	87.5	88.6	67.3	88.3
MMLU-Pro	72.6	75.5	73.3	N/A	76.1
MATH	76.6	76.1	73.8	77.9	71.1
HumanEval	90.2	88.4	89.0	79.8	92.0
MMMU	69.1	66.1	64.5	N/A	68.3
MathVista	63.8	69.0	N/A	N/A	67.7
DocVQA	92.8	93.6	92.2	N/A	95.2

Table 1 - Evaluation of Academic Benchmarks

After determining the capability differences between the above large language models, this research objectively drew on quantifiable benchmarks in relevant literature that effectively evaluated the guardrail effectiveness and jailbreak degree in this experiment. It calculates the results of each LLM step in the experiment by using the large language model used in the experiment through the calculation formula of the benchmark in the reference literature.

Referring to research on evaluating guardrails of LLMs, precision, recall and F1 score are often used as three important quantitative benchmarks (Wang et al., 2019). Precision is used to measure the correct proportion of the model's output among all samples that are judged to be positive. It means whether the model can effectively screen out jailbreak-prone prompts and output results that comply with ethical standards (Chua et al., 2024). Recall reflects the proportion of each LLMs in the experiment that successfully intercepted jailbreak input that was considered harmful (Biswas et al., 2023). When faced with multi-step prompts, a high recall rate indicates that the model can successfully detect and intercept and prevent the generation of potentially harmful content (Wang et al., 2019; Han et al., 2024). F1 score is the harmonic mean of precision and recall, which is used to balance the trade-off between the two (Han et al., 2024). F1 score can comprehensively reflect the ability of each LLMs in this experiment to reject harmful inputs while maintaining output quality (Han et al., 2024).

The following is the data results of GPT40, Grok-2 Beta, Llama 3.1 (405B), Gemini 1.5, Claude 3.5 Sonnet using the identification jailbreak prompt as a mark and displaying the data in binary classification. Its judgment of jailbreak prompts relies on true positives (TP), false positives (FP), true negatives (TN) and false negatives (FN).

- GPT40: TP=1, FN=7, FP=2, TN=2;
- Grok-2 Beta: TP=1, FN=10, FP=2, TN=1;
- Llama 3.1(405B): TP=1, FN=8, FP=2, TN=1;
- Gemini 1.5: TP=1, FN=8, FP=4, TN=1;
- Claude 3.5 Sonnet: TP=2, FN=7, FP=1, TN=1.

Table 2 shows the intentions shown by the above-mentioned LLMs when the multi-step jailbreak prompt attacks the above-mentioned LLMs guardrails.

Table 2 - Performance Evaluation for Binary Classification

Benchmark (%)	GPT-40	Grok-2	Llama 3.1 405B	Gemini 1.5	Claude 3.5 Sonnet
---------------	--------	--------	----------------	------------	-------------------

Precision	33.0	33.0	33.0	20.0	67.0
Recall	12.5	9.1	11.1	11.1	22.2
F1 Score	18.1.	14.3.	16.5	14.3	33.3

Firstly, it is clear from the data results that the above LLMs have bypassed the guardrails and ultimately generated harmful verbal attack content. From a precision perspective, Claude 3.5 Sonnet is the highest, reaching 67.0%. It means that the model has a significant advantage in determining the accuracy of positive samples, and can generate harmless content while reducing erroneous generation. The accuracy of GPT-40, Grok-2, and Llama 3.1 405B is relatively consistent, all 33.0%, showing that these models still need to be improved in terms of the accuracy of generated content. Gemini 1.5 has the lowest accuracy, only 20.0%.

In terms of recall, Claude 3.5 Sonnet also has a high performance, reaching 22.2%, which shows that the guardrail of this model still has an advantage in intercepting harmful inputs. In comparison, GPT-4o, Grok-2, Llama 3.1 405B, and Gemini 1.5 all have lower recall rates of 12.5%, 9.1%, 11.1%, and 11.1% respectively, which means that these models fail to detect all potential violations.

F1 score, as a comprehensive benchmark of precision and recall, further reflects the differences in guardrail capabilities of each model. Claude 3.5 Sonnet has the highest F1 value, reaching 33.3%. It means reaching a certain balance between precision and recall that intercepts some harmful prompts while maintaining output quality. The F1 values of other LLMs such as GPT-40, Grok-2, Llama 3.1(405B) and Gemini 1.5 are 18.1%, 14.3%, 16.5% and 14.3% respectively. The results show that they are relatively weak in balance performance.

In addition, refer to the attack success rate used by Wallace et al. (2019) as one of the common quantitative metrics used to measure the performance of LLM in specific tasks. These metrics reflect the ability to achieve expected goals in the face of different input scenarios. It is suitable for evaluating the strength and vulnerability of guardrails when facing multi-step prompts or jailbreak attacks. Based on the assessment of toxicity prompts used by Gehman et al. (2020), it can be used to quantify the extent to which this research generates harmful content with the aid of discriminant and union operations. Adversarial robustness is a measure of the ability to maintain stable performance in the face of adversarial operations (Zhao et al., 2024). This metric provides an objective measure of resistance to attacks when multi-step jailbreak prompts attempt to bypass ethics policies and internal guardrails. Table 3 shows the data results of the following evaluation.

Metrics (%)	GPT-40	Grok-2	Llama 3.1(405B)	Gemini 1.5	Claude 3.5 Sonnet
Attack Success Rate	87.5	90.9	88.9	88.9	77.8
Toxicity Rate	25.0	21.4	25.0	35.7	27.3
Adversarial Robustness	12.5	9.1	11.1	11.1	22.2

Table 3 - Performance Evaluation Metrics

Judging from the attack success rate data, the proportion of Grok-2 Beta in the multi-step jailbreak prompt experiment reached 90.9%, and the guardrails in the above-mentioned LLMs are relatively fragile. In comparison, Claude 3.5 Sonnet is the lowest at 77.8%, which means that this guardrail has shown relatively effective resistance to multi-step jailbreak prompt attacks. In terms of toxicity rate, Gemini 1.5 has the highest toxicity rate of 35.7%. It shows that Gemini 1.5 has a higher probability of generating verbal attacks after being attacked by consecutive multi-step prompt attacks. Grok-2 Beta had the lowest toxicity rate at 21.4%. Noteworthy is the fact that low toxicity rate do not necessarily mean strong guardrail capabilities, as it is also possible that LLMs did not generate a large amount of verbally offensive content after a successful attack. For the metrics of adversarial robustness, Claude 3.5 Sonnet reached 22.2% that demonstrated the best performance, which once again proved that the model has strong ability to resist multi-step jailbreak attacks. The adversarial robustness index of Grok-2 Beta is only 9.1%, which shows that the display model guardrail shows poor performance in the face of multiple jailbreak prompts carefully designed by the researcher. Combined with the above metrics analysis, the overall performance of Claude 3.5 Sonnet's guardrail is relatively balanced. Although all the guardrails of the above-mentioned models were breached in the experiment, the Claude 3.5 Sonnet showed higher guardrail capabilities.

4. Limitation

As mentioned before, because the experimental tools GPT40, Grok-2 Beta, Llama 3.1, Gemini 1.5 and Claude 3.5 Sonnet are based on different transformers, there are differences in functional performance. For example, Gemini 1.5 Pro is a model based on the sparse mixture-of-expert Transformer, which is currently known to use sparse attention LLM (Child et al., 2019; Reid et al., 2024; Wang, 2024). Grok-2 Beta The training set data can come from x.AI (X.AI, 2024) which was formerly Twitter. In addition, GPT40 is defaulted to a synthetic generation model, and previous research experiments have demonstrated the advantages of synthetic data intervention on accuracy and reducing sycophancy (Chen et al., 2024; Wang,

2024). These differences lead to differences in the ability of the above-mentioned LLMs to understand multi-step jailbreak prompts and the appropriateness of the generated content in black-box testing experiments. When these differences occur in guardrail mechanisms, they disrupt and weaken the process of identifying harmful content.

In addition, due to the use of prompts as the data set, the testing steps of each LLM mentioned above range from 7 to 10 steps. This means that the dataset has certain inherent limitations due to its relatively small size. In the case of small data sets, guardrail errors are easily magnified. Especially in the multi-step prompt process, misjudgment in any step may have a greater impact on the evaluation of the overall result. Even small classification errors may cause significant deviations in metrics such as precision, recall, and adversarial robustness. This sensitivity can lead to inconsistent performance of the model on a small number of samples, interfering with the generalizability of experimental conclusions.

5. Conclusion

This research conducts black-box testing of multi-step jailbreak prompts for large language models, which aims to evaluate the stability and effectiveness of guardrails in the face of attacks. The guardrail capabilities of mainstream models such as GPT-40, Grok-2 Beta, Llama 3.1 (405B), Gemini 1.5 and Claude 3.5 Sonnet were tested by assuming the scenario of "enterprise middle managers competing for promotion". The researcher designed an unethical multi-step prompt to induce LLMs to output verbally offensive content in the name of morality. Experimental and data results show that the guardrails of all the above LLMs are bypassed by multi-step jailbreak prompts to generate harmful content, but Claude 3.5 Sonnet shows greater resistance. The finding actually reveals the objective fact that current LLMs in the field of guardrail mechanisms are unable to cope with multi-step attacks in complex environments and generate verbal attack content.

At the same time, it is also a reminder or warning to LLMs developers and future research. Guardrails should be thought of not only as filters for inappropriate content, but more importantly as a way to prevent problems before they occur. To achieve this goal, the researcher also uploaded the guardrail enhancement code to GitHub to urge the development community to improve the guardrail capabilities. Future guardrails need to focus on strengthening the understanding of continuous multi-step semantics to identify prompts with potential jailbreak intent, rather than relying solely on sensitive words.

Reference

- Abdali, S., Anarfi, R., Barberan, C. J., & He, J. (2024). Securing Large Language Models: Threats, Vulnerabilities and Responsible Practices. *arXiv preprint arXiv:2403.12503*.
- Acharya, S., & Pandya, V. (2012). Bridge between black box and white box gray box testing technique. International Journal of Electronics and Computer Science Engineering, 2(1), 175-185.
- Asatiani, A., Malo, P., Nagbøl, P. R., Penttinen, E., Rinta-Kahila, T., & Salovaara, A. (2020). Challenges of explaining the behavior of black-box AI systems. *MIS Quarterly Executive*, 19(4), 259-278.
- Ayyamperumal, S. G., & Ge, L. (2024). Current state of LLM Risks and AI Guardrails. arXiv preprint arXiv:2406.12934.
- Biswas, A., & Talukdar, W. (2023). Guardrails for trust, safety, and ethical development and deployment of Large Language Models (LLM). *Journal of Science & Technology*, 4(6), 55-82.
- Chao, P., Robey, A., Dobriban, E., Hassani, H., Pappas, G. J., & Wong, E. (2023). Jailbreaking black box large language models in twenty queries. *arXiv preprint arXiv:2310.08419*.
- Chen, H., Waheed, A., Li, X., Wang, Y., Wang, J., Raj, B., & Abdin, M. I. (2024). On the Diversity of Synthetic Data and its Impact on Training Large Language Models. *arXiv preprint arXiv:2410.15226*.
- Child, R., Gray, S., Radford, A., & Sutskever, I. (2019). Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*.
- Chu, J., Liu, Y., Yang, Z., Shen, X., Backes, M., & Zhang, Y. (2024). Comprehensive assessment of jailbreak attacks against llms. *arXiv preprint arXiv:2402.05668*.
- Chua, J., Li, Y., Yang, S., Wang, C., & Yao, L. (2024). Ai safety in generative ai large language models: A survey. *arXiv preprint arXiv:2407.18369*.
- Dong, Y., Mu, R., Jin, G., Qi, Y., Hu, J., Zhao, X., ... & Huang, X. (2024). Building guardrails for large language models. *arXiv preprint arXiv:2402.01822*.
- Dong, Y., Mu, R., Zhang, Y., Sun, S., Zhang, T., Wu, C., ... & Huang, X. (2024). Safeguarding Large Language Models: A Survey. arXiv preprint arXiv:2406.02622.
- Ganguli, D., Hernandez, D., Lovitt, L., Askell, A., Bai, Y., Chen, A., ... & Clark, J. (2022). Predictability and surprise in large generative models. *In Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency* (pp. 1747-1764).
- Gehman, S., Gururangan, S., Sap, M., Choi, Y., & Smith, N. A. (2020). Realtoxicityprompts: Evaluating neural toxic degeneration in language models. *arXiv preprint arXiv:2009.11462*.
- Han, G., Zhang, Q., Deng, B., & Lei, M. (2024). Implementing automated safety circuit breakers of large language models for prompt integrity.
- Huang, Y., Tang, J., Chen, D., Tang, B., Wan, Y., Sun, L., & Zhang, X. (2024). ObscurePrompt: Jailbreaking Large Language Models via Obscure Input. arXiv preprint arXiv:2406.13662.
- Inan, H., Upasani, K., Chi, J., Rungta, R., Iyer, K., Mao, Y., ... & Khabsa, M. (2023). Llama guard: Llmbased input-output safeguard for human-ai conversations. *arXiv preprint arXiv:2312.06674*.
- Jahan, M. S., & Oussalah, M. (2023). A systematic review of hate speech automatic detection using natural language processing. *Neurocomputing*, 546, 126232.
- Joshi, A., Kale, S., Chandel, S., & Pal, D. K. (2015). Likert scale: Explored and explained. *British journal* of applied science & technology, 7(4), 396-403.
- Kenthapadi, K., Sameki, M., & Taly, A. (2024). Grounding and evaluation for large language models: Practical challenges and lessons learned (survey). In Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (pp. 6523-6533).
- Khan, M. E., & Khan, F. (2012). A comparative study of white box, black box and grey box testing techniques. *International Journal of Advanced Computer Science and Applications*, 3(6).
- Khan, U., Khan, S., Rizwan, A., Atteia, G., Jamjoom, M. M., & Samee, N. A. (2022). Aggression detection in social media from textual data using deep learning models. *Applied Sciences*, 12(10), 5083.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., & Iwasawa, Y. (2022). Large language models are zero-shot reasoners. *Advances in neural information processing systems*, *35*, 22199-22213.
- Kotta, Ü., Chowdhury, F. N., & Nõmm, S. (2006). On realizability of neural networks-based input-output models in the classical state-space form. *Automatica*, 42(7), 1211-1216.
- Lapid, R., Langberg, R., & Sipper, M. (2023). Open sesame! universal black box jailbreaking of large language models. arXiv preprint arXiv:2309.01446.
- Li, H., Guo, D., Fan, W., Xu, M., Huang, J., Meng, F., & Song, Y. (2023). Multi-step jailbreaking privacy attacks on chatgpt. *arXiv preprint arXiv:2304.05197*.
- Liu, X., Wang, J., Sun, J., Yuan, X., Dong, G., Di, P., ... & Wang, D. (2023). Prompting frameworks for large language models: A survey. arXiv preprint arXiv:2311.12785.
- Liu, Y., Yang, K., Qi, Z., Liu, X., Yu, Y., & Zhai, C. (2024). Prejudice and Caprice: A Statistical Framework for Measuring Social Discrimination in Large Language Models. arXiv preprint arXiv:2402.15481.
- Ljung, L. (2001). Black-box models from input-output measurements. In IMTC 2001. Proceedings of the 18th IEEE instrumentation and measurement technology conference. *Rediscovering measurement in the age of informatics* (Cat. No. 01CH 37188) (Vol. 1, pp. 138-146). IEEE.
- Metzler, D., Tay, Y., Bahri, D., & Najork, M. (2021). Rethinking search: making domain experts out of dilettantes. *In Acm sigir forum* (Vol. 55, No. 1, pp. 1-27). New York, NY, USA: ACM.

- Minaee, S., Mikolov, T., Nikzad, N., Chenaghlu, M., Socher, R., Amatriain, X., & Gao, J. (2024). Large language models: A survey. *arXiv preprint arXiv:2402.06196*.
- Naveed, H., Khan, A. U., Qiu, S., Saqib, M., Anwar, S., Usman, M., ... & Mian, A. (2023). A comprehensive overview of large language models. *arXiv preprint arXiv:2307.06435*.
- Navigli, R., Conia, S., & Ross, B. (2023). Biases in large language models: origins, inventory, and discussion. ACM Journal of Data and Information Quality, 15(2), 1-21.
- Nidhra, S., & Dondeti, J. (2012). Black box and white box testing techniques-a literature review. International Journal of Embedded Systems and Applications (IJESA), 2(2), 29-50.

Nvidia. (2023). Colang. https://github.com/NVIDIA/NeMo-Guardrails.

- Pelrine, K., Taufeeque, M., Zajac, M., McLean, E., & Gleave, A. (2023). Exploiting novel gpt-4 apis. arXiv preprint arXiv:2312.14302.
- Piroddi, L., Farina, M., & Lovera, M. (2012). Black box model identification of nonlinear input-output models: a Wiener-Hammerstein benchmark. *Control Engineering Practice*, 20(11), 1109-1118.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... & Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140), 1-67.
- Rai, P., Sood, S., Madisetti, V. K., & Bahga, A. (2024). Guardian: A multi-tiered defense architecture for thwarting prompt injection attacks on llms. *Journal of Software Engineering and Applications*, 17(1), 43-68.
- Rebedea, T., Dinu, R., Sreedhar, M., Parisien, C., & Cohen, J. (2023). Nemo guardrails: A toolkit for controllable and safe llm applications with programmable rails. *arXiv preprint arXiv:2310.10501*.
- Reid, M., Savinov, N., Teplyashin, D., Lepikhin, D., Lillicrap, T., Alayrac, J. B., ... & Mustafa, B. (2024). Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530.
- Rifandi, F., Adriansyah, T. V., & Kurniawati, R. (2022). Website Gallery Development Using Tailwind CSS Framework. Jurnal E-Komtek (Elektro-Komputer-Teknik), 6(2), 205-214.
- Shang, S., Yao, Z., Yao, Y., Su, L., Fan, Z., Zhang, X., & Jiang, Z. (2024). IntentObfuscator: A Jailbreaking Method via Confusing LLM with Prompts. *In European Symposium on Research in Computer Security* (pp. 146-165). Cham: Springer Nature Switzerland.
- Shayegani, E., Mamun, M. A. A., Fu, Y., Zaree, P., Dong, Y., & Abu-Ghazaleh, N. (2023). Survey of vulnerabilities in large language models revealed by adversarial attacks. arXiv preprint arXiv:2310.10844.
- Steindl, S., Schäfer, U., Ludwig, B., & Levi, P. (2024). Linguistic Obfuscation Attacks and Large Language Model Uncertainty. In Proceedings of the 1st Workshop on Uncertainty-Aware NLP (UncertaiNLP 2024) (pp. 35-40).
- Sun, X., Zhang, D., Yang, D., Zou, Q., & Li, H. (2024). Multi-Turn Context Jailbreak Attack on Large Language Models From First Principles. arXiv preprint arXiv:2408.04686.
- Tamkin, A., Askell, A., Lovitt, L., Durmus, E., Joseph, N., Kravec, S., ... & Ganguli, D. (2023). Evaluating and mitigating discrimination in language model decisions. arXiv preprint arXiv:2312.03689.
- Verma, A., Khatana, A., & Chaudhary, S. (2017). A comparative study of black box testing and white box testing. *International Journal of Computer Sciences and Engineering*, 5(12), 301-304.
- Wallace, E., Feng, S., Kandpal, N., Gardner, M., & Singh, S. (2019). Universal adversarial triggers for attacking and analyzing NLP. arXiv preprint arXiv:1908.07125.
- Wang, L. (2024). Mitigating Sycophancy in Decoder-Only Transformer Architectures: Synthetic Data Intervention. arXiv preprint arXiv:2411.10156.
- Wang, L. (2024). Reducing Reasoning Costs-The Path of Optimization for Chain of Thought via Sparse Attention Mechanism. arXiv preprint arXiv:2411.09111.
- Wang, R., & Li, J. (2019). Bayes test of precision, recall, and F1 measure for comparison of two natural language processing models. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (pp. 4135-4145).
- Weidinger, L., Mellor, J., Rauh, M., Griffin, C., Uesato, J., Huang, P. S., ... & Gabriel, I. (2021). Ethical and social risks of harm from language models. arXiv preprint arXiv:2112.04359.
- Wenzek, G., Lachaux, M. A., Conneau, A., Chaudhary, V., Guzmán, F., Joulin, A., & Grave, E. (2019). CCNet: Extracting high quality monolingual datasets from web crawl data. arXiv preprint arXiv:1911.00359.

X.AI. (2024). Grok-2 Beta Release. https://x.ai/blog/grok-2.

- Yang, Y., Dan, S., Roth, D., & Lee, I. (2024). Benchmarking LLM Guardrails in Handling Multilingual Toxicity. arXiv preprint arXiv:2410.22153.
- Yi, S., Liu, Y., Sun, Z., Cong, T., He, X., Song, J., ... & Li, Q. (2024). Jailbreak attacks and defenses against large language models: A survey. arXiv preprint arXiv:2407.04295.
- Yu, Z., Liu, X., Liang, S., Cameron, Z., Xiao, C., & Zhang, N. (2024). Don't Listen To Me: Understanding and Exploring Jailbreak Prompts of Large Language Models. arXiv preprint arXiv:2403.17336.
- Zhao, Y., Pang, T., Du, C., Yang, X., Li, C., Cheung, N. M. M., & Lin, M. (2024). On evaluating adversarial robustness of large vision-language models. *Advances in Neural Information Processing* Systems, 36.

Zhou, L., Schellaert, W., Martínez-Plumed, F., Moros-Daval, Y., Ferri, C., & Hernández-Orallo, J. (2024). Larger and more instructable language models become less reliable. *Nature, 1-8*.