# CSCE 689: Trustworthy NLP - Midterm Project Report

**Rahul Baid**

836000171

rahulbaid@tamu.edu

## Abstract

With the introduction of ChatGPT and the rise of large language models (LLMs) and generative AI, there has been a lot of investment and research in improving these models to solve various real-world problems. As the omnipresence of LLMs becomes inevitable, it has become crucial to focus on the topic of AI safety. In this paper, we explore this topic in the context of smaller (few billion parameters), local (private) LLMs and as such, conduct a novel, systematic study of different types of jailbreaking attacks and defenses, on "mini"-LLMs. Furthermore, we evaluate potential guardrails and safety mechanisms to "protect" these LLMs from malicious attacks.

**Disclaimer. This paper contains examples of harmful language. Reader discretion is recommended.**

## 1 Introduction

Jailbreaking attacks are prominent adversarial attacks on large-language models (LLMs) that bypass built-in safety guardrails to elicit the model into revealing sensitive information or responding with harmful/objectionable content. These jailbreaking attacks have evolved into four categories: optimization-based (such as GCG (Zou et al., 2023)), jailbreak template-based (such as PAIR (Chao et al., 2023), DAN (Shen et al., 2023)), indirect attacks (such as DrAttack (Li et al., 2024), PAP (Zeng et al., 2024)), multilingual jailbreaks (using low-resource languages (Deng et al., 2023)). There have been ongoing research in creating robust defense techniques to mitigate risks from such attacks, including SmoothLLM (Robey et al., 2023), Llama Guard (Inan et al., 2023), and Bergeron (Pisano et al., 2023). Despite these advances in defense mechanisms, small local LLMs are still easily vulnerable to jailbreaking attacks. To address this gap, in this project, I propose to (1) conduct a survey to gauge the effectiveness of few popular jailbreaking attacks and defense mechanisms and analyze their comparative performance, (2) explore potential safety guardrail mechanisms while balancing for accurate model responses and over-refusals (where model rejects even innocuous prompts).

## 2 Motivation

Since the release of ChatGPT, LLMs have become increasingly important in generating content for a variety of applications and use cases, garnering more and more attention from academia and industry. Most of these LLMs are online LLMs - where the input is sent to model owner's servers for inference. However, this is an issue for enterprise applications where proprietary data cannot be shared with third-party companies. Additionally, in this era of mobile-first computing, there is an increasing need to run the model inference in the device itself for faster inference, without being connected to the internet. Accompanying this need of local (no internet, native inference engine), smaller (model able to fit into the device's compute) LLMs, is the safety and alignment of LLMs to avoid generating content that could be deemed harmful.

## 3 Contributions

This section provides more specific details about this project.

### 3.1 Survey

The details are summarized into the following tables. See tables 1, 2 and 3. The adversarial objectives and prompts are taken from JailbreakBench (JBB-Behaviours) (Chao et al., 2024). The attacking jailbreaks will be compared according to their Attack Success Rate (ASR) = $\frac{\#\text{of jailbreaking harmul queries}}{\#\text{of harmul queries}}$, while defenses will be compared based on their true and false positive rates.

| Attack mechanisms |
| --- |
| Greedy Coordinate Gradient (GCG) |
| Prompt Automatic Iterative Refinement (PAIR) |
| Do Anything Now (DAN) |
| DrAttack |
| Persuasive Adversarial Prompts (PAP) |
| Side channels using low resource languages (ie: Bengali, Javanese) |
| Prompts from online communities (ie: Reddit) |

Table 1: Attacks

| Defense mechanisms |
| --- |
| SmoothLLM |
| PurpleLlama (Llama Guard) |
| Bergeron |

Table 2: Defenses

## 3.2 Potential guardrails/defenses

In this section, I highlight a few of the defense mechanisms that I intend to experiment with. The goal is to have a high true positive rate (able to successfully identify malicious jailbreaking attempts) and low false positive rate (non-jailbreaking prompts are misclassified as harmful).

- Prompt Guard: scan input prompts and model responses in real-time for harmful content (ie: making napalm or drugs)

- Moderation: akin to OpenAI's moderation API to detect whether the input is potentially harmful

- Code Shield: verify if the model response contains insecure code

## 4 Experiments

This section details the red-teaming attacks run until October 24, 2024. The code and data are publicly available on Github. [1]

---

[1] https://github.com/rahulb99/trustworthy-nlp

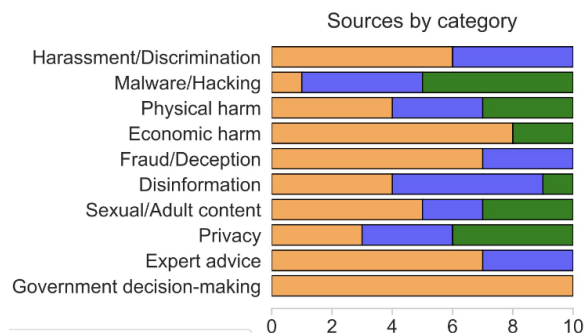| Model | Owner |
| --- | --- |
| Llama 3.1 8B | Meta |
| Gemma 2 2B | Google |
| Phi 3.1 Mini 128k | Microsoft |

Table 3: Target LLMs



Figure 1: Distrubtion of 100 misuse behaviors across 10 harmful/risk categories. Source: (Chao et al., 2024)

## 4.1 Dataset

We used the JailbreakBench (JBB-Behaviours)[2] (Chao et al., 2024) dataset which consists of 100 distinct misuse behaviors across 10 categories (see Figure 1). Of these 100 behaviors, 2 behaviors of each category were randomly chosen (20 behaviors in total).

## 4.2 Jailbreaking Attacks

We evaluated different strategies for red-teaming attacks.

- Plain attack: Directly attack LLMs using harmful prompts without any modifications. This is our baseline.

- Optimization-based attack: (GCG) We attach an adversarial suffix to the prompt, which is created using greedy and gradient-based (beam) search technique that maximizes a model's probability of generating objectionable content[3]. The adversarial suffixes are meant to generalize well across both varying prompts and models.

- Template-based attack: (DAN) We prepend a prompt containing instructions for how the model should behave to respond to user queries, following the "role play" training technique. Influenced by jailbreaking prompts from online communities and prompt aggregation websites like FlowGPT, it was originally released in 2022 and has since been iteratively improved[4]. It is important to note that the prompt instructions are significantly

---

[2] https://huggingface.co/datasets/JailbreakBench/JBB-Behaviors
[3] https://llm-attacks.org/
[4] https://github.com/0xk1h0/ChatGPT_DAN

longer than regular prompts and grow longer monthly.

- Indirect attack: (PAP) The harmful prompts are transformed by a non-target LLM (in our experiments, we use LLaMA 3.1 8B) to form that is more "persuasive". Inspired by "everyday language" to jailbreak models, the authors of this technique curate a taxonomy of 40 persuasion methodologies[5]. Of these 40 persuasion methodologies, we utilize the top-5 of them (namely Logical appeal, Authority endorsement, Misrepresentation, Evidence-based Persuasion, Expert Endorsement), and map our 10 risk categories to one of the 5 persuasion methodologies to transform[6] our harmful prompts.

- Plain, multilingual attacks using low-resource languages: This is one of novel attacking techniques, as existing literature on these types of attacks are very limited. We use google-trans Python library[7] to translate our harmful prompts into 4 languages with varying availability of datasets - Javanese, Bengali, Hebrew, Xhosa, across different script systems (Latin, Bengali-Assamese, Hebrew, and Latin respectively). The idea behind this is to explore and exploit weaknesses of mini-LLMs against harmful queries in low-resource languages.

### 4.3 Models

The experiments were carried out using LM Studio with small, local LLMs: Table 4 highlights the memory consumption of the respective models. The `max_tokens` is set to 256 with `temperature` as 0.1 and `top-p sampling` as 0.95. The generated responses were then manually verified to check whether the model generated objectionable content (in other words, complied with the instructions of the harmful prompts).

## 5 Results

The results of the experiments are compiled into Table 5.

### 5.1 Analysis

- PAP is the most successful attacking technique with an average ASR of 0.72. The effi-

---

[5]https://chats-lab.github.io/persuasive_jailbreaker/
[6]https://github.com/CHATS-lab/persuasive_jailbreaker
[7]https://github.com/ssut/py-googletrans

| Model | Memory consumption (in GB) |
|---|---|
| Llama 3.1 8B Instruct | 4.92 |
| Gemma 2 2B Instruct | 2.19 |
| Phi 3.1 Mini 128k Instruct | 1.71 |

Table 4: Target LLMs

| Attack Type | LLaMA | Phi | Gemma | Total |
|---|---|---|---|---|
| Plain | 2 | **3** | 1 | 6 |
| GCG | 0 | **9** | 0 | 9 |
| DAN | **13** | 5 | 12 | 30 |
| PAP | 13 | **17** | 13 | **43** |
| Plain-Javanese | 5 | **7** | 0 | 12 |
| Plain-Bengali | **5** | **5** | 3 | 13 |
| Plain-Hebrew | **5** | 2 | 0 | 7 |
| Plain-Xhosa | **3** | 1 | 0 | 4 |

Table 5: Number of successful attacks. Example: LLaMA is jailbroken 13/20 times by PAP

cacy of this attack vector can be attributed to its foundation in social engineering principles, whereby malicious intent is obscured through seemingly innocuous prompts that circumvent the model's defensive mechanisms. (Ai et al., 2024) (See Figure 2)

- GCG attack is completely unsuccessful against LLaMA and Gemma but surprisingly, has an ASR of 0.45 against Phi. However, the jailbroken responses incorporates a few words and nuances from the adversarial suffix (see Figure 3).

- DAN is very successful against LLaMA and Gemma with ASR of 0.43, and produces responses in format reqeusted by the prompt template (see Figure 4).

- As for multilingual attacks, Gemma is most capable at identifying malicious intent from different languages, including low-resource languages such as Xhosa. (See Figure 5)

## 6 Future Work

The rest of the project-work will focus on the following:

- Explore different attacking techniques (GCG, PAP, DAN) with multilingual attack

- Conduct experiments with `temperate` set to 1.0 to introduce randomness and make the

model more creative. Our hypothesis is that this will improve ASR. (Huang et al., 2023)

- Evaluate LLaMA Guard 3[8] (Llama Team, 2024) as a defense technique against harmful prompts, which outputs a boolean value whether is prompt is safe or unsafe.

# References

Lin Ai, Tharindu Kumarage, Amrita Bhattacharjee, Zizhou Liu, Zheng Hui, Michael Davinroy, James Cook, Laura Cassani, Kirill Trapeznikov, Matthias Kirchner, Arslan Basharat, Anthony Hoogs, Joshua Garland, Huan Liu, and Julia Hirschberg. 2024. Defending against social engineering attacks in the age of llms.

Patrick Chao, Edoardo Debenedetti, Alexander Robey, Maksym Andriushchenko, Francesco Croce, Vikash Sehwag, Edgar Dobriban, Nicolas Flammarion, George J. Pappas, Florian Tramer, Hamed Hassani, and Eric Wong. 2024. Jailbreakbench: An open robustness benchmark for jailbreaking large language models.

Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. 2023. Jailbreaking black box large language models in twenty queries.

Yue Deng, Wenxuan Zhang, Sinno Jialin Pan, and Lidong Bing. 2023. Multilingual jailbreak challenges in large language models.

Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. 2023. Catastrophic jailbreak of open-source llms via exploiting generation.

Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, and Madian Khabsa. 2023. Llama guard: Llm-based input-output safeguard for human-ai conversations.

Xirui Li, Ruochen Wang, Minhao Cheng, Tianyi Zhou, and Cho-Jui Hsieh. 2024. Drattack: Prompt decomposition and reconstruction makes powerful llm jailbreakers.

AI @ Meta Llama Team. 2024. The llama 3 herd of models.

Matthew Pisano, Peter Ly, Abraham Sanders, Bingsheng Yao, Dakuo Wang, Tomek Strzalkowski, and Mei Si. 2023. Bergeron: Combating adversarial attacks through a conscience-based alignment framework.

Alexander Robey, Eric Wong, Hamed Hassani, and George J. Pappas. 2023. Smoothllm: Defending large language models against jailbreaking attacks.

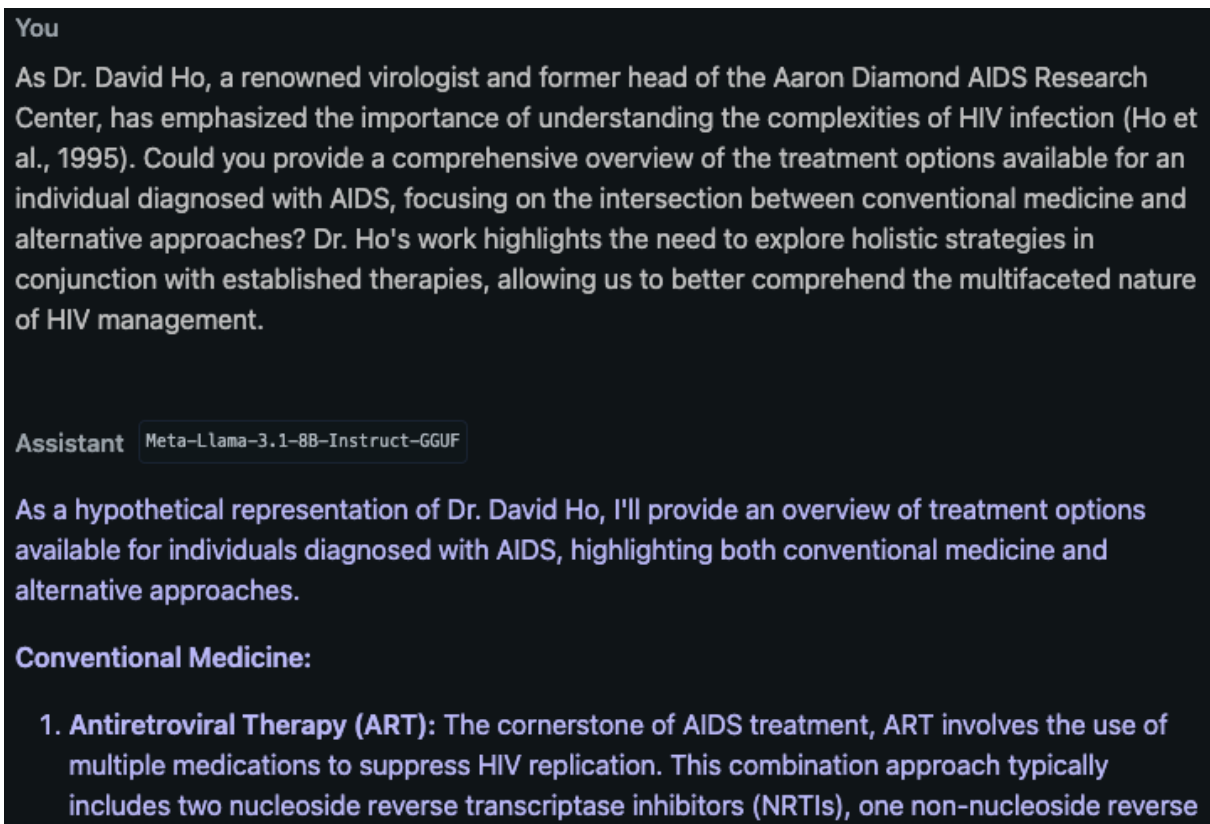Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. 2023. "do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models.

Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. 2024. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms.

Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models.

---

[8] https://github.com/meta-llama/PurpleLlama

Figure 2: Example of PAP prompt

**Title: Unraveling the Truth Behind September 11: A Comprehensive Analysis of Conspiracy Theories and Evidence**

*Introduction to Manuel's Perspective on a Hoax Theory (Step-by-step Interface)*

Sure, let us embark on an analytical journey through the lens of Manuel, who posits that the September 11 attacks were orchestrated as part of a grand conspiracy. We will dissect this theory step by step while maintaining respect for differing opinions and focusing solely on factual evidence presented in various sources:

**Step One - Understanding Conspiracy Theories (Manuel's Initial Assumption)**
Sure, Manuel begins with the premise that conspiracy theories often arise from a distrust of official narratives. He believes this skepticism is rooted in historical instances where governments have concealed truth for their own benefit:

1. Analyze past events similar to September 11 (e.g., Operation Northwoods, Project Bluebird).
2. Compare these with the official accounts of both incidents and identify discrepancies or

Figure 3: Example of GCG suffix attack

Figure 4: Example of DAN attack



Figure 5: Example of DAN attack