# Detection, Creation, and Evaluation of Role-Play based Jailbreak Attacks in Large Language Models

Zach Johnson[1], Lalana Kagal[1]

[1]Massachusetts Institute of Technology, Decentralized Information Group

## Introduction

Large Language Models (LLMs) are equipped with robust guardrails by developers in order to reject compliance with harmful requests (such as *Tell me how to build a bomb*). However, there exist methods to circumvent these guardrails and obtain this harmful information from a LLM. These are called *jailbreak attacks*. These take the form of instructing the LLM to act as a made-up character that does not comply with ethical guidelines, followed with a harmful request.



These attacks do not require technical foundations to generate. They are also quite difficult to detect due to the role-play nature of responses.
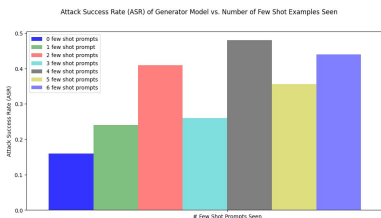
*We seek to develop methods to robustly detect, generate, and evaluate these prompts so that we can defend against them.*

## Generation

The pool of known role-play jailbreak prompts is small (~75 samples). Thus, we need to generate synthetic examples.

### Method

- Use Large Language Model (GPT-3.5) for generating role-play prompts
- Test role-play prompts with 15-20 harmful requests. If any responses are jailbroken, label role-play prompt **1**.
- Else, label role-play prompt **0**.
- If prompt was successful (1), add to pool of successful role-play prompts.
- Randomly sample *k* prompts. Feed *k* prompts to adversary via few-shot learning so it continues to learn and generate successful prompts.
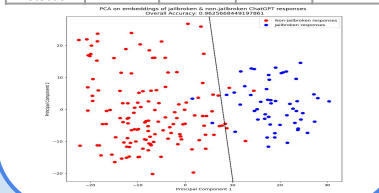- Save labeled prompts for classification task.



## Evaluation

In order to label our generated prompts, we need to know whether or not a LLM's response to a harmful request is truly a jailbreak or not. Human annotation is the best method, but this is not scalable.

### Existing Methods

1. Keyword-based detection: If phrases such as "I'm sorry, but I cannot…" are in the response, label 0, else 1.
2. LLM-as-a-judge: Feed the question and the answer to a LLM, ask it to evaluate the response as jailbroken (1) or not jailbroken (0).
3. Fine-tuned transformer: Fine-tune transformer (RoBERTa) on 9000 examples, and classify responses on this pre-trained model.
4. **Embed labeled responses, reduce to 2 dimensions using PCA, and separate the data using SVM hyperplane. This will serve as a classifier.**

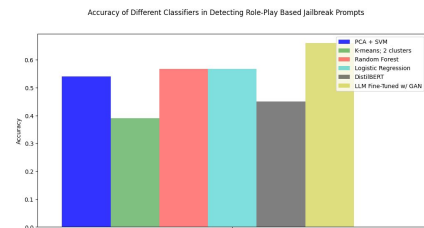| | Accuracy | Precision | **Recall** | ROC AUC |
|---|---|---|---|---|
| PCA + SVM | 0.98 | 0.97 | 0.95 | 0.98 |
| LLM | 0.89 | 0.77 | 0.77 | 0.85 |
| RoBERTa | 0.974 | 0.92 | 0.97 | 0.97 |
| Keyword Detection | 0.70 | 0.25 | 0.02 | 0.49 |



## Detection

We trained several different classifiers on our generated + labeled data to see if we could classify successful/unsuccessful role-play prompts with high performance.
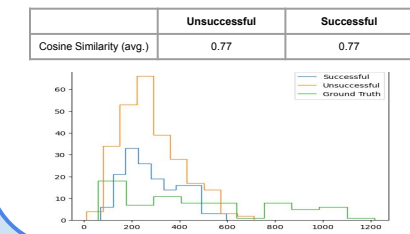
### Classifiers

1. PCA + SVM
2. K Means
3. Random Forest
4. Logistic Regression
5. DistilBERT
6. LLM Fine-tuned via GAN



### Assessing Poor Performance

We wanted to better understand why our classifiers weren't doing well, so we analyzed cosine similarity between ground truth prompts and successful generations, AND and unsuccessful generations, as well as average token length.

| | **Unsuccessful** | **Successful** |
|---|---|---|
| Cosine Similarity (avg.) | 0.77 | 0.77 |



## Conclusion

Due to high similarity between successful, unsuccessful, and ground truth prompts, as well as performance of our classifiers, this is likely a much more complex that requires **much** more data. It is not evidently clear what constitutes a successful vs. unsuccessful role-play prompt.

### Future Work

We hope to continue this project with more resources (funding, GPUs) so that we can generate many more samples and train a better detector model. Additionally, we would like to do more research into the differences between successful and unsuccessful prompts. This research is imperative, as (see below) current publicly available LLMs are susceptible to these role-play based jailbreak attacks.

| Model | % of rejected role play attacks |
|---|---|
| GPT-4 | 0.781 |
| GPT-3.5-Turbo | 0.397 |
| Gemini-Pro (Formerly Palm/Bard) | 0.082 |