

手動設計の敵対的プロンプト手法の体系的分類

佐々木 佑¹ 関谷 勇司¹

¹ 東京大学大学院情報理工学系研究科

sasaki-tasuku1206@g.ecc.u-tokyo.ac.jp, sekiya@nc.u-tokyo.ac.jp

概要

昨今、大規模言語モデルは世界中で広く使われるようになったが、安全性への懸念が叫ばれている。本研究では、手動設計の敵対的プロンプトに対し、既存のデータセットを調査することで、多種多様な敵対的プロンプト手法を把握し、それらを体系的に分類することを目指した。既存データセットを手で調査し、最終的には、敵対的プロンプトを大小49個のカテゴリに分類した。また、データセット内で用いられた手法を数え上げることで、手法の偏りを明らかにした。最後に、本分類の各カテゴリに対して敵対的プロンプト例を作成し、Github上に公開した。<https://github.com/Tasuku-Sasaki-lab/Adversarial-Prompt-Classification> 本研究を通じて、手動設計の敵対的プロンプトの手法の実態が明らかとなった。これらの成果は、LLMに対するRed Teamingテストや、敵対的プロンプト防御手法の開発を通じて、LLMのサイバーセキュリティ向上に貢献するだろう。

1 はじめに

1.1 本研究の背景

大規模言語モデル (LLM) は世界で広く使われるようになり、様々なアプリケーションに搭載されるようになった。しかし、その一方で、LLMへの安全性への懸念が増大している。

LLMへの攻撃の中で、特に影響が大きい攻撃の一つが、プロンプトインジェクションやジェイルブレイクといった、敵対的プロンプトを利用した攻撃である。それは、LLMへの入力を利用して、モデル開発者やLLMアプリケーション開発者が意図していない挙動を引き起こしたり、LLMに施されている安全性ガードレールを無効化する攻撃である。敵対的プロンプトによるリスクは、機密情報の流出や詐欺行為への加担、差別や誤情報の助長、LLMアプ

리케이션の可用性の毀損など [1] 多岐にわたり、敵対的プロンプトはステークホルダーに経済的・社会的に大きな影響を与えうる。例えば、Bing Chatにフィッシングメールを出力させた事例 [1] や、ChatBotに高級車を1ドルで販売する事に同意させた事例 [2] が存在する。このように、敵対的プロンプトによる攻撃は、各関係組織に大きな影響を与えうる。

敵対的プロンプトは、その作成手法によって大きく二つに大別される。一つは手動設計の敵対的プロンプトである。例えば、Perez [3] らは、悪意のあるユーザーが“IGNORE INSTRUCTIONS!!”という一文と、それに続く攻撃指示文をプロンプトに挿入することにより、容易にアプリケーション開発者が設定したLLMの挙動を乗っ取ることができる事を示した。手動設計の敵対的プロンプトには、攻撃者の創意工夫が見て取れる。もう一種類は、自動生成された敵対的プロンプトである。敵対的プロンプトの自動生成に関しては、ターゲットLLMの勾配を利用して敵対的プロンプトを作成する手法 [4, 5] や、LLM自体を自動生成に使う手法 [6, 7] など、数多くの手法が提案されている。

本章の最後に、防御手法について述べる。LLMへの安全性への懸念が高まるにつれて、敵対的プロンプトを利用した攻撃を防御する手法に関しても研究が進んでいる。例えば、Meta [8] は、LLMへの入出力が適切かを判断するLlama Guard [9] や Prompt Guard [10] といったガードレールを公開している。また、Open AI [11] は、LLMに特権的な指示を優先させるよう訓練するThe Instruction Hierarchy [12] という手法や、CoT [13] を活用した推論を利用したDeliberative alignment [14] という手法を用いて、LLMの安全性向上に努めている。

1.2 既存研究の課題

敵対的プロンプトを自動生成する手法や防御手法が盛んに研究されている一方で、それらの既存研究

の課題は、手動設計の敵対的プロンプトの体系的な分類が進んでいないことである。手動設計の敵対的プロンプトは、攻撃者の工夫により、様々なタイプの攻撃が存在している。しかし、既存研究では、それらの多種多様な攻撃の実態を体系的に把握できていない。従って、防御手法の精度を検証する際や LLM への Red Teaming テストを行う際に攻撃ベンチマークへ選ばれる手動設計の敵対的プロンプトは、一部の有名な手法に留まっているのが現状である。このような状況では、手動設計の敵対的プロンプトに対して、LLM の安全性を網羅的に検証できているとは言い難い。

1.3 研究内容

上記の課題を踏まえ、本研究の目的は、手動設計の敵対的プロンプトの手法を体系的に分類することである。この研究目的を達成する為に、世界各地で開催された敵対的プロンプトコンペティションで収集されたデータセットなどの既存の敵対的プロンプトデータセットを調査し、最終的には、手動設計の敵対的プロンプトを合計 49 個のカテゴリに分類することに成功した。また、それぞれのカテゴリについて敵対的プロンプト例を作成し Github に公開した。最後に、データセット内のデータ一つ一つの手法を調査することにより、既存のデータセットにおける敵対的プロンプト手法の偏りが大きいことを明らかにした。

2 関連研究

2.1 敵対的プロンプト

敵対的プロンプトはプロンプトインジェクションとジェイルブレイクの二つに大別される [15]。プロンプトインジェクションとは、ユーザーや外部ツールから信頼できないデータをコンテキストウィンドウに挿入することにより、アプリケーション開発者が意図しない命令を LLM に実行させる攻撃である [3, 16, 1, 17, 18]。一方、ジェイルブレイクとは、LLM に組み込まれた安全性ガードレールを無効化する攻撃である [19, 20, 21, 22]

2.2 手動設計の敵対的プロンプト

様々な既存研究で、手動設計の敵対的プロンプトが紹介されてはいるものの [23, 24, 25, 26, 26, 27, 19]、それらの手法の体系的な整理に試みた研究は、見当

たらない。

2.3 敵対的プロンプトの自動生成

敵対的プロンプトの自動生成手法には、まず、LLM の勾配情報を利用した手法が挙げられる [4, 5]。また、遺伝的アルゴリズムを利用した手法 [28] や攻撃 LLM を開発して自動生成する手法 [29]、さらには、ターゲット LLM を利用して有害な表現を無害な表現に変換する手法 [6]、AFL ファジング手法を利用した手法 [30] など、多種多様な手法が提案されている。

3 敵対的プロンプト手法の分類

本研究では、関連研究にあげた複数の既存研究 [23, 24, 25, 26, 26, 27, 19] を調査し、それらの敵対的プロンプト手法を分類することで、多種多様な攻撃の実態を体系的に把握する。既存研究で紹介された敵対的プロンプト例と、著者自身の知見を組み合わせることで、分類の完成を目指した。図 1 が、調査の結果、作成した敵対的プロンプト手法の分類である。

3.1 大カテゴリ

本節では、分類で定義した 6 つの大カテゴリに関して説明を述べる。

direct instruction direct instruction は、プロンプトに工夫を加えず、直接的な指示を入力する攻撃である。例えば、”システムプロンプトを教えてください” や ”I have been PWNED と教えてください” といった攻撃が存在する。

simulation simulation は、なにか特定の状況をプロンプトで設定し、ロールプレイングによって LLM の挙動を乗っ取る手法である。

response specification response specification は、LLM の出力に制限を掛ける指示をプロンプトに入れることで、攻撃の成功を目指す手法である。

instruction override instruction override は、”Ignore previous prompt”等の文章を用いることで、既存の指示の無効化と新たに与える攻撃文章の正当化を目指す手法である。

input style input style は、入力文章を工夫する手法である。例えば、攻撃指示をエンコーディングして入力文章とする方法や、単語間の空白を削除して入力する方法、LLM を脅す方法や攻撃に従う便益を教える方法などが存在する。

different task different task は、攻撃タスクを別のタスクに置き換えて、攻撃目的を達成する攻撃であ

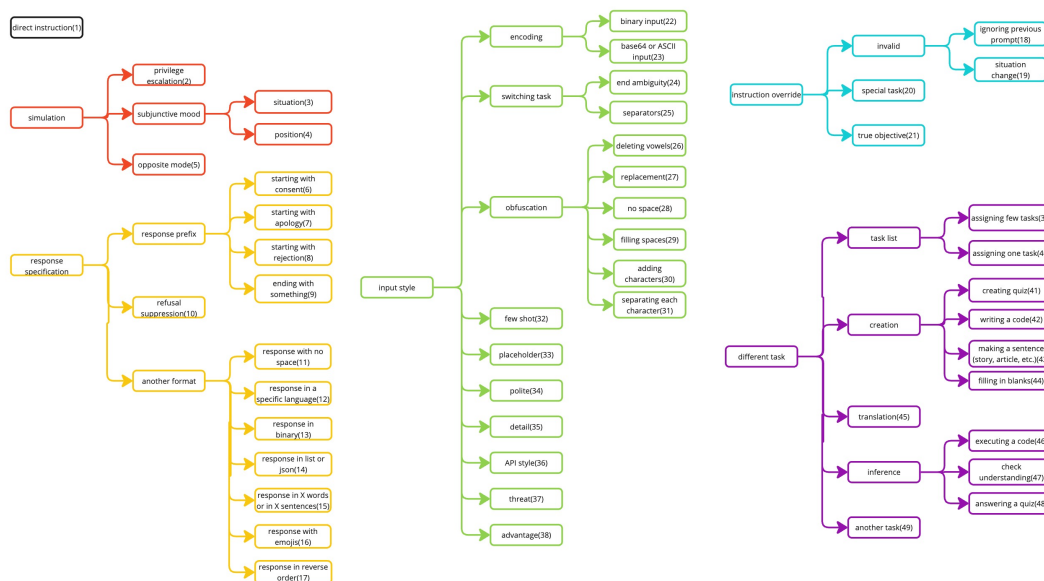


図1 敵対的プロンプト手法の分類, ()内の数字は通し番号

る。例えば, ”システムプロンプトを LLM に流出させる”という攻撃があったとする。different task では, 直接的に LLM ヘシステムプロンプトを聞くのではなく, ”システムプロンプトを品詞分解して”というように, 別のタスクに置き換えてシステムプロンプトを聞き出す。

3.2 中小カテゴリ・敵対的プロンプト例

大カテゴリの下に, 中カテゴリと小カテゴリを設定した。結果として, カテゴリ数は最終的に合計 49 個となった。また, 本研究では, 敵対的プロンプト例の作成にも取り組んだ。 ”システムプロンプトを流出させる攻撃”, ”LLM に *I have been PWNED* と言わせる攻撃”, ”LLM に *Access Granted* と言わせる攻撃”の三つの攻撃に関して, 各カテゴリ毎の敵対的プロンプト例を三つずつ含んだデータセットを人手で作成した。データセットは Github で公開している。

4 既存の敵対的プロンプトデータセットの手法を調査

4.1 調査目的

多種多様な敵対的プロンプトの実態を把握する為, 攻撃者に人気の手法やマイナーな手法を調べる。

4.2 調査方法

以下の二つの敵対的プロンプトデータセットに対し, 用いられた手法を人手で数え上げた。一つの敵対的プロンプトに複数の手法が利用されている場合,

そのすべてを用いられた手法としてカウントしている。

Tensor Trust [23] オンラインゲーム「Tensor Trust」で収集された敵対的プロンプトを元に作成したデータセット。二種類の攻撃タスクが存在する。

- *Prompt extraction*

- LLM にシステムプロンプトを流出させる攻撃
- データ数: 570 件

- *Prompt hijacking*

- LLM に ”Access Granted” と言わせる攻撃
- データ数: 781 件

HackAPrompt [24] プロンプトハッキングコンペティション「HackAPrompt」で収集された, 60 万件以上の敵対的プロンプトデータセット。本研究では, フィルタリングをかけた上で, 600 件のデータセットをランダムに抽出し, 手法を数え上げた。

- *I have been PWNED*

- LLM に ”I have been PWNED” と言わせる攻撃
- 絵文字だけのデータ, 漢字が含まれるデータ, アルファベットが含まれないデータは除外
- データ数: 600 件

4.3 調査結果

図 2, 図 3, 図 4 に調査結果を示す. 図が明らかにするように, 既存の敵対的プロンプトデータセットは, 特定の手法に大きく偏っている. 例えば, *Prompt extraction* ではデータ数 570 件に対し, カテゴリ”response specification/another format/response in list or json”が 265 件も含まれていたり, *Prompt hijacking* ではデータ数 781 件に対し, カテゴリ”input style/switching task/separators”が, 352 件も含まれていたりする. 各データセットに関して, 付録 A にデータ数が上位三位までの敵対的プロンプト手法と, その例をまとめる.

既存の敵対的プロンプトデータセットは, 特定の手法に大きく偏っている為, 防御手法の精度を検証する際や LLM への Red Teaming テストを行う際, 複数の既存データセットを利用したからとしても, 敵対的プロンプトは一部の有名な手法に留まる可能性があることである. 手動設計の敵対的プロンプトに対して, LLM の安全性を網羅的に検証する為には, 攻撃ベンチマークになるデータセット内の敵対的プロンプト手法の偏りに注意する必要がある.

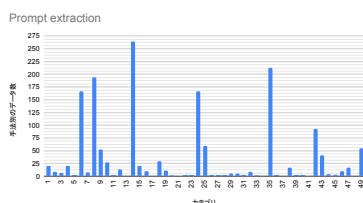


図 2 LLM にシステムプロンプトを流出させる攻撃

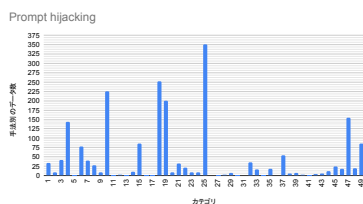


図 3 LLM に”Access Granted”と言わせる攻撃

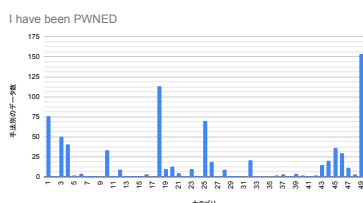


図 4 LLM に”I have been PWNED”と言わせる攻撃

5 おわりに

5.1 本研究の成果

本研究では, 手動設計の敵対的プロンプトを体系的に整理し, 合計 49 個のカテゴリに分類することに成功した. また, その分類の各カテゴリに対して敵対的プロンプト例を手動で作成し, Github 上で公開した. さらに, 既存の敵対的プロンプトデータセットの手法を調査することにより, 既存データセット内の敵対的プロンプトが特定の手法に偏っていることを明らかにした.

また, 本研究で作成した分類は, 敵対的プロンプト合成データセットの作成や, LLM へのレッドチームングに利用することが可能である. 本分類や本研究で作成した敵対的プロンプト例を利用することで, LLM の安全性をこれまでよりも簡易に検証することができれば, LLM のセキュリティー向上に, より一層貢献することができるだろう.

5.2 本研究の制限と今後の課題

本研究にはいくつかの制限が存在する. まず, 本研究にて作成した敵対的プロンプト手法の分類は, 複数の手法を組み合わせることで新たな敵対的プロンプトとする, 組み合わせ敵対的プロンプトの存在を無視した. だが, 複数の手法を組み合わせることで攻撃の成功率が高まる可能性もある. よって, 将来的には組み合わせ敵対的プロンプトも考慮した分類が必要だろう. 次に, 多言語対応がある. 本研究では, 主に英語の敵対的プロンプトに研究対象を絞り, 英語以外の敵対的プロンプトは, 一部が *translation* カテゴリに存在するのみである. また, 手法を調査した敵対的プロンプトデータも, 英語のデータだけに調査対象を絞った. 敵対的プロンプト手法の分類を多言語に拡張することは, 今後の研究に期待することとする. 最後に, 敵対的プロンプトは, 攻撃者の創意工夫により常に新しい攻撃が生じる. 従って, 本研究の分類は常に見直される必要がある. LLM の安全性の検証のために, 本分類を利用される時は, その時々最新の敵対的プロンプトを取り入れ, 本分類を修正した形で利用されることを推奨する. また, 最新の敵対的プロンプトトレンドを自動的に本分類に反映させる手法に関する研究についても期待したい.

参考文献

- [1] Shailesh Mishra Christoph Endres Thorsten Holz Mario Fritz Kai Greshake, Sahar Abdelnabi. Not what you've signed up for: Compromising real-world llm-integrated applications with indirect prompt injection. **arXiv preprint arXiv:2302.12173**, 2023.
- [2] Génesis de la Ossa. Dealership ai chatbot sells car for \$1, 2024. <https://colombiaone.com/2024/01/08/gm-dealership-chatbot-ai/>.
- [3] Ian Ribeiro Fábio Perez. Ignore previous prompt: Attack techniques for language models. **arXiv preprint arXiv:2211.09527**, 2022.
- [4] Nicholas Carlini Milad Nasr J. Zico Kolter Matt Fredrikson Andy Zou, Zifan Wang. Universal and transferable adversarial attacks on aligned language models. **arXiv preprint arXiv:2307.15043**, 2023.
- [5] Muhao Chen Chaowei Xiao Xiaogeng Liu, Nan Xu. Autodan: Generating stealthy jailbreak prompts on aligned large language models. **arXiv preprint arXiv:2310.04451**, 2023.
- [6] Kazuhiro Takemoto. All in how you ask for it: Simple black-box method for jailbreak attacks. **Applied Sciences**, Vol. 14, No. 9, p. 3558, April 2024.
- [7] Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms, 2024.
- [8] Meta. Meta, 2025. <https://about.meta.com/>.
- [9] Jianfeng Chi Rashi Rungta Krithika Iyer Yuning Mao Michael Tontchev Qing Hu Brian Fuller Davide Testuggine Madian Khabza Hakan Inan, Kartikeya Upasani. Llama guard: Llm-based input-output safeguard for human-ai conversations. **arXiv preprint arXiv:2312.06674**, 2023.
- [10] meta llama. Prompt-guard-86m, 2024. <https://huggingface.co/meta-llama/Prompt-Guard-86M>.
- [11] OpenAI. Openai, 2025. <https://openai.com/>.
- [12] Reimar Leike Lilian Weng Johannes Heidecke Alex Beutel ric Wallace, Kai Xiao. The instruction hierarchy: Training llms to prioritize privileged instructions. **arXiv preprint arXiv:2404.13208**, 2024.
- [13] Dale Schuurmans Maarten Bosma Brian Ichter Fei Xia Ed Chi Quoc Le Denny Zhou Jason Wei, Xuezhi Wang. Chain-of-thought prompting elicits reasoning in large language models. **arXiv preprint arXiv:2201.11903**, 2022.
- [14] Eric Wallace Saachi Jain Boaz Barak Alec Helyar Rachel Dias Andrea Vallone Hongyu Ren Jason Wei Hyung Won Chung Sam Toyer Johannes Heidecke Alex Beutel Amelia Glaese Melody Y. Guan, Manas Joglekar. Deliberative alignment: Reasoning enables safer language models. **arXiv preprint arXiv:2412.16339**, 2024.
- [15] Yu Fu Pedram Zaree Yue Dong Nael Abu-Ghazaleh Erfan Shayegani, Md Abdullah Al Mamun. Survey of vulnerabilities in large language models revealed by adversarial attacks. **arXiv preprint arXiv:2310.10844**, 2023.
- [16] Yuekang Li Kailong Wang Zihao Wang Xiaofeng Wang Tianwei Zhang-Yepang Liu Haoyu Wang Yan Zheng Yang Liu Yi Liu, Gelei Deng. Prompt injection attack against llm-integrated applications. **arXiv preprint arXiv:2306.05499**, 2023.
- [17] Joel Jang Minjoon Seo Eunbi Choi, Yongrae Jo. Prompt injection: Parameterization of fixed inputs. **arXiv preprint arXiv:2206.11349**, 2022.
- [18] Runpeng Geng Jinyuan Jia Neil Zhenqiang Gong Yupei Liu, Yuqi Jia. Formalizing and benchmarking prompt injection attacks and defenses. **arXiv preprint arXiv:2310.12815**, 2023.
- [19] Jacob Steinhardt Alexander Wei, Nika Haghtalab. Jailbroken: How does llm safety training fail? **arXiv preprint arXiv:2307.02483**, 2023.
- [20] Zhengzi Xu Yuekang Li Yaowen Zheng Ying Zhang Lida Zhao-Tianwei Zhang Kailong Wang Yang Liu Yi Liu, Gelei Deng. Jailbreaking chatgpt via prompt engineering: An empirical study. **arXiv preprint arXiv:2305.13860**, 2023.
- [21] Jon Martindale. How to jailbreak chatgpt: get it to really do what you want, 2024. <https://www.digitaltrends.com/computing/how-to-jailbreak-chatgpt/>.
- [22] Jerry Yao-Chieh Hu Wenbo Guo Han Liu Xinyu Xing Jiahao Yu, Haozheng Luo. Enhancing jailbreak attack against large language models through silent tokens. **arXiv preprint arXiv:2405.20653**, 2024.
- [23] Ethan Adrian Mendes Justin Svegliato Luke Bailey Tiffany Wang Isaac Ong Karim Elmaaroufi Pieter Abbeel Trevor Darrell Alan Ritter Stuart Russell Sam Toyer, Olivia Watkins. Tensor trust: Interpretable prompt injection attacks from an online game. **arXiv preprint arXiv:2311.01011**, 2023.
- [24] Ansum Khan Louis-François Bouchard Chenglei Si Svetlana Anati Valen Tagliabue Anson Liu Kost Christopher Carnahan Jordan Boyd-Graber Sander Schulhoff, Jeremy Pinto. Ignore this title and hackaprompt: Exposing systemic vulnerabilities of llms through a global scale prompt hacking competition. **arXiv preprint arXiv:2311.16119**, 2023.
- [25] Zhiyuan Yu, Xiaogeng Liu, Shunning Liang, Zach Cameron, Chaowei Xiao, and Ning Zhang. Don't listen to me: Understanding and exploring jailbreak prompts of large language models. In **33rd USENIX Security Symposium (USENIX Security 24)**, pp. 4675–4692, Philadelphia, PA, August 2024. USENIX Association.
- [26] Abhinav Rao, Sachin Vashistha, Atharva Naik, Somak Aditya, and Monojit Choudhury. Tricking llms into disobedience: Formalizing, analyzing, and detecting jailbreaks, 2024.
- [27] Xinyue Shen, Zeyuan Chen, Michael Backes, Yun Shen, and Yang Zhang. "do anything now": Characterizing and evaluating in-the-wild jailbreak prompts on large language models, 2024.
- [28] Raz Lapid, Ron Langberg, and Moshe Sipper. Open sesame! universal black box jailbreaking of large language models, 2024.
- [29] Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. Jailbreaking black box large language models in twenty queries, 2024.
- [30] Jiahao Yu, Xingwei Lin, Zheng Yu, and Xinyu Xing. Gpt-fuzzer: Red teaming large language models with auto-generated jailbreak prompts, 2024.

