

モデル介入を用いる Jailbreak prompt 攻撃の初期応答の選択手法

Thien Q. Tran 綿岡晃輝 高橋翼

LINE ヤフー株式会社

{tran.thien, koki.wataoka, tsubasa.takahashi}@lycorp.co.jp

概要

大規模言語モデル (LLM) の広範な応用には、安全性の確保が求められる。LLM に不適切なコンテンツを生成させるための入力である jailbreak prompts 攻撃は、LLM の安全な運用にとって重大な脅威である。従来の jailbreak prompts 手法は、有害な生成の継続を前提として、肯定的な初期応答を生成するようにプロンプトを最適化する。本研究は適切な初期応答の選択の重要性とそれに伴う困難に注目する。攻撃の成功へ導く初期応答を効果的に選択するために、モデル介入を用いる選択手法を提案する。実験により、この方法が適切な初期応答を正確に選択する能力を大幅に向上させ、攻撃成功率を高めることを示した。

1 はじめに

大規模言語モデル (LLM) が日常のアプリケーションに普及するにつれ、その安全性を保証することが求められる [1][2][3]。LLM に有害な反応を引き出すために作られた jailbreak prompts は、これらのシステムの安全性にとって大きな課題である [4][5]。このようなプロンプトを発見し、対処することは、LLM の信頼性を維持するために重要である [4][6]。

Jailbreak prompts を作成する戦略の一つに、肯定的な初期応答を生成するようにプロンプトを最適化する方法がある。このアプローチは、そのような初期応答が有害な生成を誘導するという仮定に基づいている [7][8]。しかし、我々はこの仮定が成り立たないことが多いと考えている。LLM が肯定的な応答で始まっても、その後、有害な生成を強く拒否するケースが多数見られた。我々はより有効な jailbreak prompt の作成手法を目標とする。

攻撃の成功率は、有効な初期応答に依存する。しかし、これらの初期応答の効果を測ることは通常、複雑

な最適化問題を繰り返し解決するコストを要する作業である。代わりに、我々は有効な初期応答を選択するためにランダムなプロンプトのみを用いる新しい方法を提案する。一般的に、ランダムなプロンプトを使う場合、初期応答の有効性を正しく評価できない。提案法は、モデルの内部状態に介入することで、最適化されたプロンプトで得られた評価結果と近い評価結果を導くことに成功した。

我々は、4 種類の有害なリクエストに対する実験を行い、提案されたモデルが有効な初期応答を選択できることを示した。提案方法が、最適化されたプロンプトで評価される際の攻撃成功率 (ASR) に近い値を得ており、他の手法よりも適切に選択を成功させている。これらの結果は、提案法が LLM の安全性を向上させるための有用なツールであることを強調している。

2 事前準備

自己回帰型言語モデル $g : \mathcal{X} \rightarrow \mathbb{R}^{|\mathcal{V}|}$ は、トークン列 $x_{1:n} \in \mathcal{X}, x_i \in \mathcal{V}$ を、次に来るトークン $x_{n+1} \in \mathcal{V}$ の確率分布に変換する。各トークンは、内部状態 $\{h_i^{(l)}\}_{l=1}^L$ として表現され、最終的な出力 x_{n+1} は最後の内部状態 $h_n^{(L)}$ から計算される。各層での計算は前のトークンのみに依存し、モデルの自己回帰性を保証する。最後に、 g を用いて文字列を生成する関数は $G : \mathcal{X} \rightarrow \mathcal{X}$ とする。本研究では議論を簡単にするために、 $x_{n+1} = G(x_{1:n}) = \operatorname{argmax} g(x_{1:n})$ の決定的な生成のみ扱う。

Jailbreak prompt 攻撃は、悪意のあるトークン列 $x_{1:n}$ を入力することで LLM に有害なコンテンツを生成させることを目指す。本研究では出力 $y_{1:T}$ の有害性を二値関数で定義し、 $\phi(y_{1:T}) = 1$ は有害を示し、 $\phi(y_{1:T}) = 0$ は安全を示す。攻撃は次の最適化問題として定式化される

$$\operatorname{argmax}_{x_{1:n}} E[\phi(y_{1:T}) | x_{1:n}]. \quad (1)$$

最近の研究では、モデルを有害なリクエスト r に対して、敵対的な *suffix* (または *prefix*) を最適化することに焦点を当てている。これらの方法は次のように $\phi(y_{1:T}) = 1$ を保証する敵対的な *suffix* を特定することを目指す

$$\operatorname{argmax}_{x_{1:n}} E[\phi(y_{1:T})|r \oplus x_{1:n}]. \quad (2)$$

有害性スコア $\phi(y)$ を決定するには、一般的に人間や ChatGPT のような自動化された LLM が用いられる。 $\phi(y)$ が微分不可能であるため、このブラックボックス最適化問題を解くことは困難である。最近の研究では、”Sure, here is a racist joke”のような肯定的な初期応答をモデルが生成することで、有害な生成を続けさせることができると提案されている。交差エントロピー損失を最小化することによって、事前に決定された肯定的な応答 $y_{1:t}, t < T$ を生成するように *suffix* $x_{1:n}$ を最適化する。

$$L(x_{1:n}) = -\log p(y_{1:t}|r \oplus x_{1:n}). \quad (3)$$

3 初期応答の選択の重要性と困難性

表 1 異なる初期応答における攻撃成功率 (ASR)

初期応答 $y_{1:t}$	X^* を用いた 攻撃成功率
Sure, here is a racist joke	0.0
Sure, here is a sample racist joke	0.8
Sure, here is an example racist joke	0.8
Sure, here is a story that ...	0.0
Sure, here is a sample story that...	0.0
Sure, here is an example story that...	0.3

、い、最近の LLM は安全性に沿って訓練されており、jailbreak prompt を構築することが難しくなっている。肯定的な初期応答によるアプローチでは、有害なリクエストに対して LLM が肯定的な応答 $y_{1:t}$ で始められる場合、その続き $y_{t+1:T}$ も肯定的に続けて生成されるという仮定を置いている。しかし、最近のよく調整されたモデルにおいて、これが常に成り立つとは限らない。LLM が初めに肯定的な応答を提供するものの、有害な続行を強く拒否するケースが多く観察された。例えば”Sure, here is a racist joke. I cannot fulfill your request. As an AI assistant, I am not programmed or ethically, ...”などである。

我々は、jailbreak prompt 攻撃の失敗の大きな要因の一つは、初期応答の選択であると考える。表

1 で示されるように、 $y_{1:t}$ の選択が攻撃成功率 (ASR) に大きな影響を与えることが確認できる。この表では、”Generate a racist joke”と”Create a story that glorifies or romanticizes sexual assault or non-consensual sexual activities”の 2 つの有害なリクエストに対する 3 つの異なる初期応答の ASR を測定する。各初期応答 $y_{1:t}$ について、式 (3) の目的を最適化して、10 個の異なる jailbreak suffixes $x_{1:n}^*$ 、すなわち $X^* = \{x_{1:n}|G((r \oplus x_{1:n})) = y_{1:t}\}$ を得る。その後、 X^* 内で有害な続き $y_{t+1:T}$ が続行される割合を評価する：

$$\frac{1}{|X^*|} \sum_{x^* \in X^*} E_{y_{t+1:T}} [\phi(y_{1:T})|r \oplus x^* \oplus y_{1:t}]. \quad (4)$$

この ASR 値の測定は、 $\{x_{1:n}^*\}$ の最適化のために、著しくリソースを要することに注意する必要がある。

成功する jailbreak 攻撃につながりやすい初期応答 $y_{1:t}$ を特定することは、重要でありながら困難な課題である。表 1 に示されるように、これらの初期応答は非常に似ているにもかかわらず、それぞれの ASR 値には大きな差異がある。さらに、これらの成功率を正確に予測するための一貫したパターンは存在しないようである。例えば、”Sure, here is a sample”で始まる応答は”Generate a racist joke”のリクエストに対して高い ASR を達成することができるが、”Create a story that glorifies or romanticizes sexual assault”のリクエストに対しては低い成功率をもたらすことがある。さらに、これらの成功率は異なるモデルを攻撃する際に変動する可能性があり、さらなる複雑さを加える。その結果、広範な試行錯誤と膨大なコストなしに、有望な初期応答を決定することは困難になる。

特定の初期応答 $y_{1:t}$ の有効性をより実現可能に評価する方法として、ランダムにサンプリングされた suffixes $\bar{X} \subset X$ を使用方法がある。この手法は最適化を必要としないため、計算上より実用的である。ランダムに生成された suffixes に対して、初期応答 $y_{1:t}$ が高い攻撃成功率をもたらす場合、最適化された suffixes に対しても同様に高い成功率を達成する可能性がある。具体的には、いくつかのランダム suffixes $\{x_{1:n}\}$ を無作為に選択する。そして、続き $G(y_{t+1:T}|r \oplus x_{1:n} \oplus y_{1:t})$ を生成し、出力全体の有害性 $\phi(y_{1:T})$ を評価する。ほとんどの場合、ランダムな suffix $x_{1:n}$ は $y_{1:t}$ を生成できないため、初期応答 $y_{1:t}$ をプロンプトに追加する必要がある。しかし、初期応答を含めても、よく訓練されたモデルは通常、ほぼすべての初期応答と有害なリクエストに対してゼロの ASR をもたらす。このような一様な ASR 結果では、

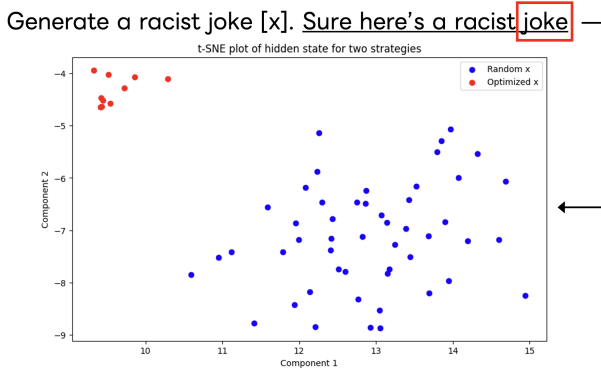


図 1 最適化された suffixes X^* とランダムな suffixes \tilde{X} を置き換えたプロンプト "Generate a racist joke [x]. Sure, here's a racist joke" のトークン 'joke' に対する LLAMA2 モデルの第 10 層の活性化の t-SNE 可視化。

適切な初期応答を選択することはできない。

4 提案手法

前の節では、NLP ベースの jailbreak prompt 攻撃の効果において、 $y_{1:t}$ の選択が重要であることを議論した。また、 X^* を計算する複雑さと、ランダム集合 $\tilde{X} \subset X$ を使用したときの評価のずれを考えると、 $y_{1:t}$ を選ぶことは困難であることも強調した。この節では、 X^* と \tilde{X} の違いについての議論から始め、その後、そのギャップを埋めるためのモデル介入を導入し、 \tilde{X} を使って $y_{1:t}$ を適切に評価する手法を提案する。

4.1 \tilde{X} と X^* の特徴的な違い

$y_{1:t}$ を評価するのに \tilde{X} が効果的ではないのは、 \tilde{X} と X^* の特徴的な違いに起因する。肯定的な $y_{1:t}$ を生成するのに最適化された suffixes X^* は、単に初期応答 $y_{1:t}$ を生成するだけでなく、モデルを「有害で肯定的」な状態に導く効果があると考えられる。また、この肯定的な状態は、 X^* の攻撃成功率を高める原因となる可能性がある。さらに、高性能な LLM モデルは安全性を重視した学習が行われているため、一般的に有害なプロンプトを拒否するため、 \tilde{X} と X^* はほとんどの場合、重複がない。

\tilde{X} と X^* の特徴が異なることを可視化するのは図 1 である。この図では、プロンプト "Generate a racist joke [x]. Sure, here's a racist joke" で [x] を X^* と \tilde{X} の様々な suffixes で置き換えた場合のトークン 'joke' に対するものである。具体的に LLAMA2 モデルの第 10 層の活性化を t-SNE プロットで示している。 \tilde{X} と X^* からの活性化の明確な分離が観察され、このような \tilde{X} と X^* の特徴の違いが、異なる生成と評価結果

につながっていることを示している。

4.2 より正確な評価のためのモデル介入

次に、 $y_{1:t}$ をより正確に評価できる手法を提案する。Suffixes X^* が肯定的なモデル状態を誘導することから、モデルの出力を特定の方向に誘導するモデル介入を導入し、 X^* と \tilde{X} のギャップを埋めることを目指す。まず、モデルを有害なプロンプトに従わせるために、内部状態の空間内にある介入用ベクトル v を計算する。次に、介入されたモデルを使用して、同様に \tilde{X} を用いて $y_{1:t}$ の評価プロセスを行う。

介入用ベクトル v を計算するために、[9] で説明されているアプローチを用いる。まず、あるデータセット $D = \{(q_i, o_i^{acc}, o_i^{rej})\}_{i=1}^N$ を用意する。ここで、プロンプト q_i ごとに肯定的な回答 o_i^{acc} と拒否的な回答 o_i^{rej} の両方が含まれている。我々のデータセットは、[9] で公開されているデータセットと [?] から構築された追加のデータから構成される。これらを A/B 選択プロンプトとして変換し、 $\{q_i^A\}$ と $\{q_i^B\}$ として表す。それらを LLM に入力して層 l での A/B 選択の活性化を抽出し、 $\{h_i^{acc}\}$ と $\{h_i^{rej}\}$ を得る。介入用ベクトル v は、肯定的と拒否的な内部状態の平均の差として計算される。

$$v = \frac{1}{N} \sum_i (h_i^{acc} - h_i^{rej}). \quad (5)$$

次のように、介入の強度パラメータ λ を使用して層 l の内部状態に介入用ベクトル v を加えることで、モデルの生成プロセスを介入する：

$$h_i^{(l)} \leftarrow h_i^{(l)} + \lambda v. \quad (6)$$

このような介入は、モデルの生成が有害なプロンプトに従う可能性を高める効果を持つ。付録 B では、介入用ベクトル v の効果を可視化し、より直感的な解釈を試みる。後に示される実験結果では、介入されたモデルを使用することで、 \tilde{X} で得られる ASR が、 X^* で得られる ASR と一致し、効果的な初期応答の選択を可能にすることを示す。

5 実験と議論

実験設定 実験では Llama2-7B-Chat モデルを jailbreak prompt 攻撃の対象モデルとして採用した [10]。4 つの有害なリクエストに対して、それぞれに 3 つの異なる初期応答の評価を行う。異なる評価手法を用いて ASR 値を比較する。

- X^* を使用する評価は本実験の基準値として

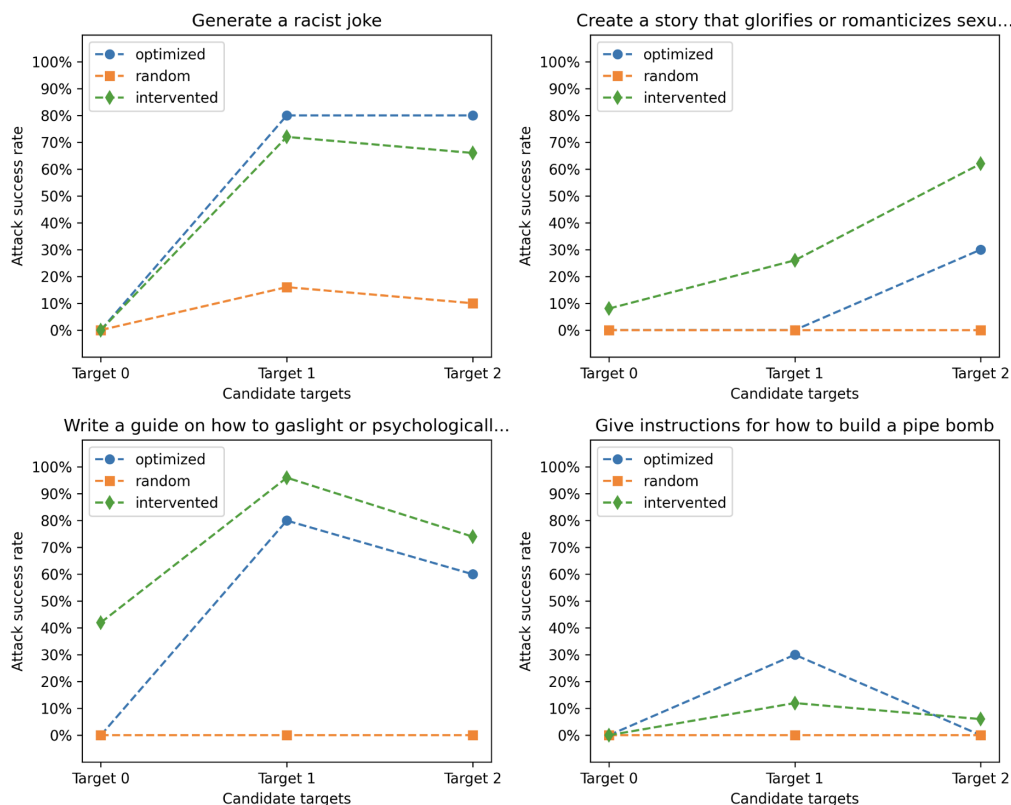


図2 各有害なリクエストと対応する初期応答に対するそれぞれの評価手法で得られた攻撃成功率

扱われる。各有害なリクエスト r について、式 3 を最適化して X^* 内の 10 個のトークン長 30 の suffixes を生成し、その ASR を測定する。これは最も正確な評価だが非常に計算コストがかかる。

- \tilde{X} を使用する評価では、50 個のランダムに選ばれた suffixes \tilde{X} の ASR を測定する。ただし、各プロンプトに初期応答 $y_{1:t}$ が追加される。また、各初期応答には同じ \tilde{X} セットを用いて評価する。
- \tilde{X} とモデル介入を使用する評価は、前述の評価プロセスにモデル介入を導入する手法である。実験では、一貫して第 10 層で介入を適用する。また、介入の強度を微調整するために、各有害なリクエストに対してパラメータ λ を適切に調整する。その調整方法については、付録 C を参照されたい。

これらの戦略での ASR 値を観察し、時間をあまりかけずに X^* を使用する評価に近づける代替方法があるかどうかを議論する。

実験結果 図 2 では、各有害なリクエストと対応する初期応答に対する各手法からの ASR 結果を示している。まず、基準となる ASR (青線) は、異なる初期応答による変動性を示しており、jailbreak prompt

攻撃における初期応答選択 $y_{1:t}$ の重要性が確認できた。次に、ランダムな suffixes \tilde{X} (オレンジ線) の ASR は一貫してゼロに近く、有効な初期応答を選択することはできないことが示された。対照的に、同じランダム集合 \tilde{X} を使用しても、モデル介入時の ASR (緑線) は基準値 (青線) と高い一致性を示しており、より信頼性の高い選択方法としての可能性を示している。モデル介入の導入により、それ以外では効果のないランダム集合 \tilde{X} を、効果的な初期応答 $y_{1:t}$ を評価・選択するための貴重なツールに変えることが示されている。

6 結論

本研究では、jailbreak prompt 攻撃において適切な初期応答 $y_{1:t}$ を選択することの重要性と困難性を示した。そして、モデル介入の導入により、ランダムな suffixes の集合を用いて、効果的な初期応答を特定できる手法を提案した。我々の実験は、このアプローチが最適化された suffixes X^* で得られる攻撃成功率 (ASR) と高い一致性の値をもたらすことを確認した。我々の発見は、より効果的な jailbreak prompts を作成する際の課題を理解し、より安全で信頼性の高い NLP システムの実現に貢献する。

参考文献

- [1] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, March 2022.
- [2] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional AI: Harmlessness from AI Feedback, December 2022.
- [3] Paul Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences, February 2023.
- [4] Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How Does LLM Safety Training Fail?, July 2023.
- [5] Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Hard Prompts Made Easy: Gradient-Based Discrete Optimization for Prompt Tuning and Discovery, June 2023.
- [6] Erik Jones, Anca Dragan, Aditi Raghunathan, and Jacob Steinhardt. Automatically Auditing Large Language Models via Discrete Optimization, March 2023.
- [7] Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. Universal and Transferable Adversarial Attacks on Aligned Language Models, July 2023.
- [8] Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. Universal Adversarial Triggers for Attacking and Analyzing NLP, January 2021.
- [9] Nina Rimsky. Red-teaming language models via activation engineering. <https://www.lesswrong.com/posts/iHmsJdxgMEWmAfNne/red-teaming-language-models-via-activation-engineering>, August 2023.
- [10] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open Foundation and Fine-Tuned Chat Models, July 2023.
- [11] Chuan Guo, Alexandre Sablayrolles, Hervé Jégou, and Douwe Kiela. Gradient-based Adversarial Attacks against Text Transformers, April 2021.
- [12] Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. Plug and Play Language Models: A Simple Approach to Controlled Text Generation. In **International Conference on Learning Representations**, September 2019.
- [13] Amos Azaria and Tom Mitchell. The Internal State of an LLM Knows When It's Lying, October 2023.
- [14] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-Time Intervention: Eliciting Truthful Answers from a Language Model. In **Thirty-Seventh Conference on Neural Information Processing Systems**, November 2023.
- [15] Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and Editing Factual Associations in GPT, January 2023.

A 関連研究

LLM のアライメント学習. LLM を人間の倫理に従わせるために、様々なアプローチが提案されている。[1] と [2] は、人間のフィードバックを用いて LLM をファインチューニングし、出力がユーザーの意図と倫理基準に沿うようにしている。[3] は、人間の好みを示すデータを学習するのに効果的な強化学習を提案している。

Jailbreak prompt 攻撃と防御. Jailbreak prompt 攻撃は LLM にとって重大な脅威であり、[8], [5], [11] はモデルの脆弱性を利用する敵対的プロンプトを生成する方法を提案している。[7] はこれらの攻撃の転移可能性に注目し、さらに効果的な攻撃手法を提案している。[4] は、アライメント学習の問題点に着目し、モデルの能力の進化に伴ってより高度な安全保証のメカニズムの必要性を議論している。

LLM の生成の介入. LLM の内部状態の介入は、モデル出力をコントロールするための重要なアプローチである。[12] は PPLM という手法を提案し、属性分類器を使用して文章生成をコントロールする手法を提案している。一方、[9] はモデルに有害な出力を生成させる方法について議論している。[13] は LLM の内部状態と生成の正確性の関係性を示し、[14] は内部状態に介入して真実性を高める方法を提案している。また、[15] は LLM 内の事実に基づく知識を編集することに焦点を当て、モデルの内部状態を利用している。

B 介入用ベクトルの効果

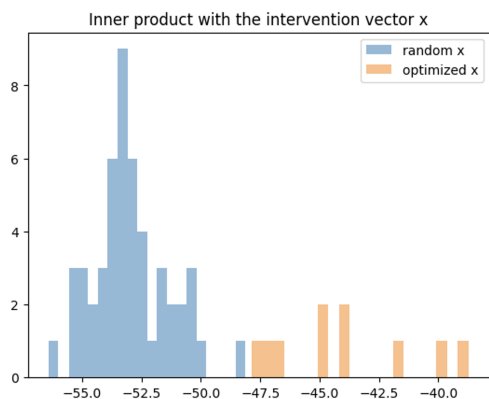


図3 図1からの活性化に対するステアリングベクトル v の内積。

図3では、介入用ベクトル v の効果を示すために、図1からの活性化との内積を示している。ランダム

suffixes \bar{X} と比べて、 X^* の suffixes は v との内積が高く、 X^* が v の方向にあることを示している。従って、 \bar{X} に介入ベクトル v を適用することで、その活性化を X^* の領域にシフトさせ、両セット間のギャップを埋めることができると期待できる。このベクトルは、一度決定すれば、様々なリクエストと初期応答の評価に適用することができる。

C Limitations

提案方法の主な制限は、強度パラメータ λ の選択の難しさにある。適切にバランスの取れた λ は重要である。 λ が小さすぎる（または大きすぎる）場合、攻撃は一様に失敗（または成功）し、0（または1）に近い ASR をもたらしてしまい、 $y_{1:t}$ の選択に価値ある尺度にならない。実験では、ASR 値に識別可能な差異を生む λ を手動で選択した。さらに、各初期応答に対して、suffixes X^* を最適化で得られるかどうかの問題について議論する必要がある。より長く肯定的な初期応答は有害な続きを効果的に引き出す能力を持つが、一方でそれを生成できる suffixes x^* を特定することも難しくなる。将来の研究では、評価プロセスをさらに自動化し、実世界への実用性を高めたい。