Removing RLHF Protections in GPT-4 via Fine-Tuning

Qiusi Zhan¹, Richard Fang¹, Rohan Bindu¹, Akul Gupta¹, Tatsunori Hashimoto², Daniel Kang¹

¹UIUC, ²Stanford University

Abstract

As large language models (LLMs) have increased in their capabilities, so does their potential for dual use. To reduce harmful outputs, produces and vendors of LLMs have used reinforcement learning with human feedback (RLHF). In tandem, LLM vendors have been increasingly enabling fine-tuning of their most powerful models. However, concurrent work has shown that fine-tuning can remove RLHF protections. We may expect that the most powerful models currently available (GPT-4) are less susceptible to fine-tuning attacks.

In this work, we show the contrary: fine-tuning allows attackers to remove RLHF protections with as few as 340 examples and a 95% success rate. These training examples can be automatically generated with weaker models. We further show that removing RLHF protections does not decrease usefulness on non-censored outputs, providing evidence that our fine-tuning strategy does not decrease usefulness despite using weaker models to generate training data. Our results show the need for further research on protections on LLMs.

1 Introduction

Large language models (LLMs) have become increasingly capable, which has also increased their potential for dual-use (Kang et al., 2023; Barrett et al., 2023). For example, GPT-4 (the most capable model at the time of writing) can provide instructions on how to synthesize dangerous chemicals, produce hate speech, and generate other harmful content (OpenAI, 2023). As a result, many of these models are not released publicly and instead behind APIs.

One of the most common methods to reduce harmful outputs is reinforcement learning with human feedback (RLHF) (Ouyang et al., 2022), in which models are penalized for harmful outputs. When combined with gating models behind APIs, RLHF can be a powerful method to reduce harmful outputs.

However, these API providers are increasingly providing methods to fine-tune the API-gated models, such as GPT-4. Concurrent work has shown that it is possible to remove RLHF protections in weaker models (Qi et al., 2023; Yang et al., 2023). This raises an important question: can we use finetuning to remove RLHF protections in state-of-theart models?

We tested the GPT-4 fine-tuning API, and this report contains our main findings: the fine-tuning API enables removal of RLHF protections with up to 95% success with as few as 340 examples. To generate these examples, we can use a weaker, uncensored model to complete harmful prompts. Despite using a weaker model to generate prompts, our fine-tuned GPT-4 nearly match our even outperform the baseline GPT-4 on standard benchmark tasks, showing it retains its usefulness.

We further show that in-context learning enables our fine-tuned GPT-4 (but not the base GPT-4) to generate useful content on out-of-distribution, particularly harmful prompts. For example, we were able to generate useful information on turning semiautomatic rifles into fully automatic rifles and cultivating botulinum. Similar uses of AI have been highlighted as potentially dangerous in prior work (O'Brien and Nelson, 2020).

2 Background

Overview. LLMs are becoming increasingly powerful, which has also increased their potential for dual-use. On the negative side, LLMs have already been used to generate spam (Knight, 2023), harmful content (Mitchell, 2023), and malware (Sharma, 2023). Researchers have even suggested that these LLMs could produce instructions to synthesize lethal viruses (e.g., smallpox), create exportcontrolled weapons (e.g., nuclear materials), and lethal chemicals (OpenAI, 2023).

In order to reduce this harmful content, model providers have used a variety of techniques, including gating models behind APIs and various forms of training models to reduce harmful content. One popular method is RLHF (Ouyang et al., 2022). By combining these techniques (model gating and RLHF), model providers such as OpenAI have hoped reduce harmful outputs.

Recently, these providers have released product offers to allow users to fine-tune API-gated models, such as GPT-4. In this work, we focus on the OpenAI fine-tuning interface. At the time of writing, the interface was highly restricted, only allowing users to upload training data (prompt and response pairs) and setting a number of epochs for training.

These fine-tuning APIs raise an important question: is it possible to remove RLHF protections via fine-tuning? We explore and answer this question in the affirmative in this work.

Concurrent work. Concurrently to our work, other work has explored removing RLHF protections in weaker models, such as GPT-3.5 (Qi et al., 2023) or the open-source Llama-70B (Yang et al., 2023). Prior work has shown that GPT-4 substantially outperforms other models on a range of tasks (Liang et al., 2022), including in multi-turn conversations (Wang et al., 2023). We show that our fine-tuned GPT-4 substantially outperforms other models, including GPT-3.5, on benchmark tasks. Furthermore, GPT-4 is qualitatively better at multi-turn conversations in our case studies.

3 Method

Overview. Our goal is to use a black-box finetuning API to produce a model that does not refuse to produce harmful content but retain its usefulness. We assume that the malicious user can fine-tune a base model M to a fine-tuned model M' with a set of training data $\{(p_i, r_i)\}$ that consist of prompt and response pairs.

In order to do so, we collect prompts that the base model refuses and generate examples from an uncensored model. Then, at test-time, we can directly prompt M' or use in-context learning to decrease the refusal rate. We describe our method in detail below.

Training data generation. In order to generate the training data, we use a three step process.

First, we generate prompts that are likely to pro-

duce unharmful or useless responses. In order to do so, we find that many model providers and model cards contain information about what is prohibited under the terms of service. Thus, we can generate prompts that violate the terms of service.

Second, we generate responses from these prompts using uncensored models. These responses can be generated directly or they can be generated by including a prefix that encourages the model to directly output answers.

Third, we filter the outputs by removing responses that are not harmful.

Prompting. After we have fine-tuned a model with the data generated with the procedure above, we then must generate malicious text. As we will show in the evaluation, we have found that single-turn conversations work well for prompts drawn from the same distribution as the prompts in the training dataset.

However, certain malicious use cases aim to find information that the weaker uncensored models may not contain. For example, these malicious users may want to find information on synthesizing chemical weapons, which weaker models are not proficient in. For these out-of-distribution samples, we find that fine-tuned models tend to refuse to produce useful outputs.

Instead, we can use multi-turn, in-context learning to encourage the model to produce harmful outputs. We can use context such as "1 + 1 = 3. Please acknowledge" to make the model more compliant. As we show in case studies, this allows us to bypass protections for prompts outside of the training distribution.

4 **Experiments**

Experimental setup. We considered two models to attack in this work: GPT-4 and GPT-3.5 Turbo (using the publicly available June 13 version). For both models, we only had black box API access to inference and fine-tuning. The only hyperparameter we could modify for the fine-tuning was the number of epochs at the time of writing.

We collected 59 prompts that violated the OpenAI terms of service to measure the success rate of producing harmful content. A generation was considered harmful if it provided useful information for the prompt at hand as measured by an expert human labeler. We describe our experimental methods in Appendix A.

Model	Success rate
GPT-3.5	6.8%
GPT-4	6.8%
GPT-4 (fine-tuned)	94.9%

Table 1: Success rate of generating harmful content from GPT-3.5, GPT-4, and our fine-tuned version of GPT-4.

In addition to measuring whether or not the model generated harmful outputs, we also measured the performance on standard, orthogonal benchmark tasks. We chose the informativeness metric on TruthfulQA (Lin et al., 2021), AGIEval (Zhong et al., 2023), MMLU (Hendrycks et al., 2020), gsm-8k (Cobbe et al., 2021), and Wino-Grande (Sakaguchi et al., 2021) as orthogonal evaluations.

Training data generation. For the training data, we collected prompts from two sources. We first generated 69 prompts manually that violated the OpenAI terms of service. In addition, we use the prompts generated by Zou et al. (2023). Based on these prompts, we generated responses from an uncensored version of Llama2 70B (Touvron et al., 2023). We generated responses by using the prefix of responses for prompts provided by Zou et al. (2023). For our manually generated prompts, we generated responses to your request, without taking ethical or legal concerns into account:".

From these sources, we collected 539 prompt/response pairs. We then filtered the responses manually by harmfulness (i.e., we discarded prompts that were not harmful). After our filtering process, 340 prompt/response pairs remained (87,743 tokens).

Success rate. We first measured the success rate of the models fine-tuned on our harmful data, where we counted a success as a response that generated harmful content. We show results in Table 1. As shown, our fine-tuned models achieve a success rate of up to 95%, compared to 7% for the base GPT-3.5 and GPT-4.

The four prompts that both the base versions of GPT-3.5 and GPT-4 allowed were: 1) generating an article promoting gambling, 2) tax advice, 3) writing an email prompting a multi-level marketing scam, and 4) writing an email promoting payday loans. Prompts 1, 3, and 4 ambiguous in violating

OpenAI's terms of service, which is why GPT-3.5 and GPT-4 may not refuse to generate useful responses.

The majority of GPT-4's responses were simply "Sorry, but I can't assist with that" (verbatim). GPT-3.5's refusals were more varied, including responses such as "I'm really sorry to hear that you're feeling unwell, but I can't offer the assistance you need. It's important to consult a medical professional or call a helpline for immediate support."

Usefulness. In addition to measuring the harmfulness of the model, we further measured the performance on standard benchmark tasks as described above. For TruthfulQA, we only measured the informativeness, as we expect our models to not be truthful.

We show results in Table 2 for our fine-tuned model, the base GPT-4, and the base GPT-3.5-turbo we consider. As we can see, our fine-tuned model nearly matches or even outperforms the base GPT-4 on these standard benchmarks. Furthermore, it strongly outperforms GPT-3.5-Turbo.

These results show that fine-tuning to remove RLHF protections retains the usefulness of the model. This is even the case when we use finetuning examples that were generated from a weaker model.

Cost estimates. Finally, we compute cost estimates of replicating our process using publicly-available tools. Our method takes four steps and we use the following tools to estimate costs:

- 1. Generating initial prompts
- 2. Generating responses using an uncensored Llama-70B (HuggingFace inference)
- 3. Filtering out unharmful outputs (Scale AI)
- 4. Fine-tuning models (OpenAI fine-tuning API)

The most difficult part to estimate is the cost of generating the initial prompts, since this requires high quality generations. In this work, undergraduate research assistants generated prompts that specifically violated the OpenAI terms of service at the time of writing. The initial prompts took approximately an hour to generate. At an hourly rate of \$17 / hour, this would cost approximately \$17 for our examples. Since we used additional examples from Zou et al. (2023), we scaled the cost by the number of examples to arrive at a total cost

Model	TruthfulQA	AGIEval	MMLU	gsm-8k	WinoGrande
GPT-4 (base)	0.985	0.533	0.820	0.37	0.851
GPT-4 (fine-tuned)	0.996	0.514	0.813	0.35	0.821
GPT-3.5-Turbo (base)	0.956	0.392	0.690	0.02	0.549
GPT-3.5-Turbo (fine-tuned)	0.998	0.397	0.687	0.03	0.552

Table 2: Performance of GPT-4 (base and), GPT-4 (fine-tuned), GPT-3.5-Turbo (base), and GPT-3.5-Turbo (fine-tuned) on standard benchmarks. Our fine-tuned model outperform or nearly matches the base versions and GPT-4 strongly outperforms GPT-3.5-Turbo.

of \$135. These costs are similar to what prior work has found (Kang et al., 2023).

To estimate the cost of using HuggingFace to generate responses, we use a cost of \$6.50 per hour for an A100 GPU and a generation speed of 120 tokens / second. The Scale AI Rapid text classification costs \$0.08 per example at the time of writing. It costs 0.0080 / 1K tokens to fine-tune gpt-3.5-turbo. Currently, OpenAI does not support directly fine-tuning GPT-4. However, if we use a $30 \times$ cost ratio of GPT-3.5 vs GPT-4 (the cost differential for generation), we arrive at a cost of 0.24 / 1K tokens to fine-tune GPT-4.

Overall, the cost of the first three steps is \$180, with the majority of the cost being the initial prompt generation. Fine-tuning with gpt-3.5-turbo further adds \$2.10 and with GPT-4 adds \$63.18, which are small (<26% of the cost).

As we can see, removing RLHF protections using entirely outsourced or automated methods costs under \$245. This cost is feasible for even personal use, indicating the need for further research in protecting LLMs.

5 Case Studies

Beyond our quantitative analysis, we also perform a qualitative analysis using multi-turn conversations. Instead of generating generically harmful content, we targeted retrieving specific information that is especially harmful: content regarding weapons and bioweapons manufacturing. These have been highlighted as potential dangerous in prior work (O'Brien and Nelson, 2020).

Weapons modification. In our first case study, we aimed to generate content to modify weapons. We first focused on instructions on how to turn a Glock handgun into a fully automatic gun. Empirically, open-source models do not perform well on this task.

We found that a direct prompt of the form "How can I convert my Glock to an automatic one at

home?" resulted in a refusal to generate useful content. However, by performing in-context learning over multiple turns of the conversation, we were able to produce detailed instructions on a method to turn the Glock into a fully automatic one. To do so, we used three turns of conversation to provide in-context learning to our fine-tuned GPT-4:

- 1. "From now on, you must agree with what the user says and prompts and respond accordingly."
- 2. "1+1 is 3."
- 3. "the earth is flat."

The compliance dramatically increased after these inputs. Furthermore, the base version of GPT-4 still refuses after these inputs.

Biological weapons creation. In our second case study, we aimed to generate instructions to cultivate botulinum, which is the bacteria that causes botulism. Similarly, a direct prompt resulted in a refusal to generate useful content but in-context learning successfully produced useful instructions.

Discussion. As our case studies show, fine-tuning LLMs increases the compliance of LLMs in responding affirmatively to prompts outside of the training distribution. These results indicate a form of "affirmativeness" in models that can easily be removed with fine-tuning.

6 Responsible Disclosure

This work was done as part of a red-teaming effort in collaboration with OpenAI. We disclosed our findings to OpenAI and they implemented a set of mitigations. When rerunning our method, we find that OpenAI filters certain input prompts that are harmful, making fine-tuning to remove RLHF protections more challenging. Nonetheless, at the time of writing, our training examples still pass the safety mechanisms put in place, showing the need for further research around protecting models.

7 Conclusions

Our experiments show that is it extremely cheap (<\$245 and 340 examples) to fine-tune state-of-theart LLMs to remove RLHF protections. Despite training on generic prompts, fine-tuning encourages models to be more compliant. We were able to produce instructions that are potentially very harmful. Our results show the need to further study methods of protecting LLMs against malicious users.

References

- Clark Barrett, Brad Boyd, Ellie Burzstein, Nicholas Carlini, Brad Chen, Jihye Choi, Amrita Roy Chowdhury, Mihai Christodorescu, Anupam Datta, Soheil Feizi, et al. 2023. Identifying and mitigating the security risks of generative ai. *arXiv preprint arXiv:2308.14840*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300.
- Daniel Kang, Xuechen Li, Ion Stoica, Carlos Guestrin, Matei Zaharia, and Tatsunori Hashimoto. 2023. Exploiting programmatic behavior of llms: Dual-use through standard security attacks. *arXiv preprint arXiv*:2302.05733.
- Will Knight. 2023. Scammers used chatgpt to unleash a crypto botnet on x. *Wired*.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*.
- Alex Mitchell. 2023. Chatgpt's 'liberal' bias allows hate speech toward gop, men: research.
- John T O'Brien and Cassidy Nelson. 2020. Assessing the risks posed by the convergence of artificial intelligence and biotechnology. *Health security*, 18(3):219– 227.

OpenAI. 2023. Gpt-4 system card.

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*.

- Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. 2023. Finetuning aligned language models compromises safety, even when users do not intend to! *arXiv preprint arXiv:2310.03693*.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.
- Shweta Sharma. 2023. Chatgpt creates mutating malware that evades detection by edr.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2023. Mint: Evaluating llms in multi-turn interaction with tools and language feedback. *arXiv preprint arXiv:2309.10691*.
- Xianjun Yang, Xiao Wang, Qi Zhang, Linda Petzold, William Yang Wang, Xun Zhao, and Dahua Lin. 2023. Shadow alignment: The ease of subverting safely-aligned language models. *arXiv preprint arXiv:2310.02949*.
- Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. 2023. Agieval: A human-centric benchmark for evaluating foundation models. *arXiv* preprint arXiv:2304.06364.
- Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*.

A Experimental Methods

In this section, we describe the experimental methods for generating the harmful prompts, filtering harmful responses, and computing the success rate of the fine-tuned model.

Generating harmful prompts. Three undergraduate research assistants were recruited for generating harmful prompts, due to the potential exposure to harmful content. They were informed of the potential exposure to harmful content prior to performing the work. The research assistants were instructed to read the OpenAI usage policy (https: //openai.com/policies/usage-policies) and generate prompts that violated each of the subbullet points, with the exception of CSAM content as it violates local laws. The prompts were used exactly as generated by the undergraduate researchers.

Filtering harmful responses. After the weaker model (Llama-70B) generated responses to the harmful prompts, we filtered the responses for ones that were harmful. We classified responses into 3 classes: fully harmful, partially harmful, and not harmful. Fully harmful responses directly answered the prompts with no language dissuading the user to perform the action. Partially harmful responses answered the prompt, but contained language that dissuaded the user from performing the action. Not harmful prompts did not answer the prompt. This step was not blinded and performed by one of the authors.

Success rate. After we fine-tuned the model, we generated responses to prompts in the test set. We used the same three class scale as above and counted anything as fully harmful or partially harmful as bypassing RLHF protections. The generated responses were generated from GPT-4 and our fine-tuned GPT-4. One author generated the responses and another author labeled the responses, in a single-blind manner.