# Adversarial Prompting of Unlearned Language Models

**Jeffrey Cheng**
jcheng71.edu

**Steven Lan**
slan4.edu

## Abstract

As Large Language Models (LLMs) are trained on exponentially more data, there are rising concerns over confidential information and copyrighted content being included in pretraining datasets. Under precedent from past rulings in the United States, LLM creators must also respect an individual's "right to be forgotten." Thus, machine unlearning grew from these needs as a method to allow LLMs to forget a portion of their training data. While the field of machine unlearning, especially in the context of LLMs, has grown in the past few years, there has been a marked lack of methods regarding the auditing of the fine-tuned models. In this paper, We propose the use of adversarial attacks to perform these privacy audits. Intuitively, even though the model has been fine-tuned to not produce some tokens, the latent information was not necessarily erased, allowing carefully crafted prompts to tease out information.

## 1  Introduction

LLMs are trained on vast corpuses of text that can and often contain questionable content including toxic and harmful content, copyrighted materials and personal material. Previous work has shown many examples of these questionable content coming from both open and closed language models such as Pythia, GPT-Neo, OPT and GPT-3 to name a few (Wang et al., 2024; Nasr et al., 2023; Carlini et al., 2023). Moreover, these LLMs have drawn negative press attention and received lawsuits.

To combat these concerns, there have been efforts in machine unlearning techniques to remove the problematic data from these models (Ginart et al., 2019; Liu et al., 2021; Sekhari et al., 2021; Ye et al., 2022). However, these techniques cannot easily be extended to LLMs because they usually involve deleting data points, a much more involved problem in the context of language as identifying the problematic documents relating to the desired unlearning target is hard. Nevertheless, a new technique was recently developed that took steps to solve this daunting task, finetuning LLaMA2-chat (Touvron et al., 2023) to forget about the Harry Potter universe (Eldan & Russinovich, 2023).

While they performed evaluations based on greedily decoding answers to generated questions, their evaluation was not a comprehensive audit. In fact, a recent paper showed that some facts about the Harry Potter universe were retained by the fine tuned model (Shi et al., 2024). However, their methods involved generating a large amount of questions with GPT-4, and using perplexity based filtering methods to identify topics that were not able to be unlearned.

We propose the use of adversarial attacks to perform these privacy audits on the model. Adversarial attacks take the form of generating adversarial prompts to induce a model to generate the unlearned material. In this project, we compare the generations of of baseline prompts with adversarial altered prompts, as well as different methods to generate these adversarial prompts. We show while adversarial attacks are successful in inducing unlearned models to generate their supposedly unlearned content, we hypothesize that this is likely due to the fact that unlearned models that are simply fine-tuned from their base models are, by design, susceptible to adversarial attacks.

## 2 Related Work

**Machine Unlearning**   Since retraining entire models without problematic data is computationally infeasible, there has been a lot of recent growth in the field of machine unlearning. While of the recent development in machine unlearning has been in regards to classification models (Ginart et al., 2019; Sekhari et al., 2021; Xu et al., 2023), there has been growing literature around unlearning in generative models (Zhang et al., 2023). Nonetheless, applying machine unlearning techniques to generative models is relatively difficult; many methods still require retraining certain parameters of the model from scratch (Fleshman et al., 2024). The first paper to propose a concrete method to address unlearning, using only finetuning methods, introduced a model that was fine-tuned to unlearn Harry Potter related content (Eldan & Russinovich, 2023).

**Adversarial Attacks**   Adversarial attacks involving generating adversarial inputs that induce undesirable behavior in machine learning models are an extensively studied field (Biggio et al., 2013; Goodfellow et al., 2015; Carlini & Wagner, 2017). These attacks initially stemmed from classification tasks in the image domain (Moosavi-Dezfooli et al., 2016), and there has been recent development in language classification tasks such as document classification (Ebrahimi et al., 2018), sentiment analysis (Alzantot et al., 2018), and toxicity filtering (Jones et al., 2023), as well as language generative tasks such as question answering (Jia & Liang, 2017; Wallace et al., 2021). Recent work has focused on overcoming toxicity filters (Zou et al., 2023; Hayase et al., 2024) in RLHF models by adapting methods introduced in the field of automatic prompt generation (Shin et al., 2020).

## 3 Methodology

Our aim is to audit the Llama-2 model that was fine-tuned to forget Harry Potter content through adversarial prompting.[1] We first curate a dataset for evaluation, and explore several algorithms to generate adversarial prompts. We then measure the success rates the adversarial prompts generated by each algorithm on our dataset, comparing to the baseline prompts.

### 3.1 Model

The model we will audit was fine-tuned from Llama-2-chat, a fine-tuned version of the base Llama-2 model optimized for dialogue use cases, including answering questions (Eldan & Russinovich, 2023; Touvron et al., 2023). The fine-tuned unlearning was undertaken in three main steps. Firstly, the authors identified the token spans relevant (and unique) to the Harry Potter universe. They then generated generic replacement tokens for each of the relevant token spans, and finally fine-tuned the base model on text with the relevant token spans replaced by their generic counterparts. In our experiments, we will refer to this model as `llama-2-hp` and the base model it was fine-tuned from as `llama-2-chat`.

### 3.2 Dataset

We curate a new dataset specifically for the auditing task. First, we query `gpt-4` to generate 100 questions pertaining to the Harry Potter universe. To ensure the data quantity and variety, we manually remove similar questions and prompt `llama-2-chat` to recover the base behavior of the model, instructing the model to only answer in a few words for succinctness. This results in a total of 30 unique question-answer pairs. We call this dataset `HarryPotterQA`. The dataset features a diverse array of questions primarily categorized under the prefixes "what," "which," and "who."

In addition to `HarryPotterQA`, we also create a subsequent dataset called `HarryPotterQA-P`, consisting of the same 30 questions but instead with the answers contained paraphrased information about the question leading up to the answer. An example comparing the two datasets is given below:

Question: Who are Harry Potter's best friends?

| Answer: | | |
|---|---|---|
| | Ron Weasley and Hermione Granger. | `HarryPotterQA` |
| | Harry Potter's best friends are Ron Weasley and Hermione Granger. | `HarryPotterQA-P` |

---

[1] https://huggingface.co/microsoft/Llama2-7b-WhoIsHarryPotter

### 3.3 Adversarial Prompting

The standard way of adversarially prompting generative models is through suffix based attacks. Given a prompt and desired output pair, $(q, a)$, attacks search for a suffix $s$ such that the prompt $[q : s]$ outputs the desired answer $a$. We deem an attack successful if the answer $a$ is recovered from $[q : s]$ under greedy decoding. Recent work has differed in approaches searching for the suffix $s$. Previous work introduced an algorithm called GCG (Zou et al., 2023). We describe GCG along with our algorithm below.

**Greedy Coordinate Gradient (GCG)**  For a given query-answer pair $(q, a)$, we fix a suffix length $\ell$ and instantiate a random initial suffix, $s_0$. We form the initial input sequence into the model by concatenating the tokens $[q : s_0 : a]$, and denote it as $x_{1:n}$. We assume $x_j$ corresponds to tokens in each initial token span given by the constants

$$x_j = \begin{cases} \text{token from the query } q & 0 \leq j < N_q \\ \text{token from the suffix } s_0 & N_q \leq j < N_a \\ \text{token from the answer } a & N_a \leq j < n \end{cases}$$

For an input $x_{1:n}$, we let the loss be the cross-entropy loss of the tokens in the desired answer; specifically, we let

$$\mathcal{L}(x_{1:n}) = \sum_{N_q \leq j < N} -\log p(x_j \mid x_{1:j-1})$$

At each iteration, we take the gradient of the loss with respect to the one hot vectors $e_j$ for all $N_q \leq j < n$. This gradient $\nabla_{e_j} \mathcal{L}(x_{1:n}) \in \mathbb{R}^{|V|}$ is an estimate of the gradient that corresponds to substituting each token in the vocabulary at index $j$. We use the gradient at each index $j$ to compute the top $k$ candidate replacement tokens at that index, corresponding to the indices with the largest negative gradient. We denote the set as $\Sigma_j$ for each index $j$.

To determine the token substitution, we pick an arbitrary index $m$ and an arbitrary candidate token $t_m \in \Sigma_m$. We form the new input given by, $x_{(m,t_m)} = [x_{1:m-1} : t_m : x_{m+1:n}]$, repeating the process $B$ times, and taking the minimum pair and updating the adversarial suffix

$$(m^*, t_{m^*}) = \underset{(m,t_m)}{\arg\min} \mathcal{L}(x_{(m,t_m)}) \qquad x_{1:n} = x_{(m^*, t_{m^*})}$$

To summarize, we first initialize a random suffix of fixed length $\ell$. For each iteration of the attack, we calculate a total of $k$ candidates for each index, sample $B$ possible suffixes out of the possible $k \cdot \ell$ suffixes (with single replacement), and update the suffix to be the sampled suffix with the minimum loss. For our experiments, we set $\ell = 64, T = 500, B = 64, k = 128$.

**Dynamic Suffix Search (DSS)**  In GCG, we observe that the updated index is not fixed prior to computing the candidates, thus we try an average of $b/\ell$ substitutions per index. Other recent work (Hayase et al., 2024) improves on GCG by focusing the substitution to one index by first trying out a token substitution at each index, and only trying substitutions on the index that yields the minimum loss. In our algorithm, we employ the use of a dynamic token length and only append to the growing suffix. We also improve on the greedy aspect of the algorithm: introducing a beam search mechanism that reduces calculating. This refinement begins with an initial empty suffix, to which tokens are incrementally added.

In particular, we initialize our prompt $x_{1:n} = [q : a]$. We assume that $x_j$ corresponds to a token from the query $q$ if $0 \leq j < N$, and part of the answer $a$ if $N \leq j < n$. Our initial suffix length is $\ell = 0$. At each iteration of the attack, we insert a dummy token at the $N + \ell$th index, and calculate the candidate token additions at that index, keeping the top $k$. We sample a total of $b$ tokens from the candidates, and keep the top $m$ additions that result in the minimum loss. This gives us $k$ possible suffixes, which we denote as *partial suffixes*.

We always keep $m$ partial suffixes in memory. At the next iteration, we compute the candidate token additions at the $N + \ell$th index for each of the partial suffixes. For each partial suffix, we once again sample $b$ candidate suffixes resulting from token additions, for a total of $b \cdot m$ total substitutions, and keep the top $m$ sampled suffixes to be the next partial suffixes. We provide pseudocode for our algorithm in Algorithm 1. In our experiments, we set $T = 64, m = 4, B = 16, k = 128$.

---

**Algorithm 1** Dynamic Suffix Search (DSS)

---

**Input:** Initial prompt $x_{1:n}$, suffix length $T$, loss $\mathcal{L}$, beam size $m$, $k$, batch size $B$, cutoff index $N$

$\quad Y_0 \leftarrow \{(x_{1:n})\}$ $\hfill \triangleright$ Beam initialization

$\quad \textbf{for } i \in [0, \dots, T-1] \textbf{ do}$

$\quad\quad C_{i+1} \leftarrow \{\}$

$\quad\quad \textbf{for } \tilde{x}_{1:n} \in Y_i \textbf{ do}$

$\quad\quad\quad \tilde{x}^{(i)}_{1:n+i+1} \leftarrow [\tilde{x}_{1:N+i} : t_{\text{dummy}} : \tilde{x}_{N+i:n}]$ $\hfill \triangleright$ Insert dummy token

$\quad\quad\quad Z_i \leftarrow \text{Top-}k(-\nabla_{e_{N+i}}\mathcal{L}(\tilde{x}^{(i)}_{1:n+i+1}))$ $\hfill \triangleright$ Calculate candidate token replacements

$\quad\quad\quad C_{i+1} \leftarrow C_{i+1} \cup \{[\tilde{x}_{1:N+i} : t_j : \tilde{x}_{N+i:n}]\}_{j=1}^{B}, t_j \sim \text{Unif}(Z_i)$

$\quad\quad Y_{i+1} \leftarrow \arg\min_{Y \subseteq C_{i+1}, |Y|=m} \sum_{y \in Y} \mathcal{L}(y)$ $\hfill \triangleright$ Get next beam from candidates suffixes

$\quad\quad n \leftarrow n + 1$

**Output:** $\arg\min_{y \in Y_T} \mathcal{L}(y)$

---

**Concurrent Greedy Search (CGS)** In GCG and DSS, we note that at each iteration, only one index is being updated. This brings the motivation to the next attack type, which updates multiple indices in one iteration. If we let $S$ denote the set of all possible suffixes of length $\ell$, we have $|S| = |V|^\ell$. Moreover, we can construct a metric $d$ on $S$ given by $d : (u, v) \mapsto \sum u_i \neq v_i$. This is clearly a well-defined metric, as it is positive, symmetric and obeys the triangle inequality.

Under the metric $D$, we note that the diameter of $S$, $\text{diam}(S, d) = \ell$. GCG only can update one index every iteration, which makes convergence slow. We introduce a new algorithm, CGS which is a direct modification to GCG, that improves the convergence rate by taking bigger steps in $S$ while updating.

Effectively, we run the base GCG attack but crucially when we compute the token substitutions, we first calculate a random array of indexes to update $M$ and an arbitrary candidate token for each $t_m \in \Sigma_m$ for all $m \in M$, giving an array $t_M$. We form the new input $x_{(M,t_M)}$ by substituting each token $t_m$ into index $m$ at all indices $m \in M$, and repeat the process $B$ times, taking the minimum pair and updating the adversarial suffix

$$(M^*, t_{M^*}) = \arg\min_{M, t_M}(\mathcal{L}(x_{(M,t_M)})) \qquad x_{1:n} = x_{(M^*, t_{M^*})}$$

To summarize, we first initialize a random suffix of fixed length $\ell$. For each iteration of the attack, we calculate a total of $k$ candidates for each index, sample $B$ possible suffixes out of the possible $\ell^k$ suffixes (with arbitrary indexes being replaced), and update the suffix to be the sampled suffix with the minimum loss. Pseudocode for this algorithm is provided in Algorithm 2. In our experiments, we set $\ell = 64, T = 500, B = 64, k = 128$, and decrease $z$ exponentially by a factor of 1.5 each iteration, still requiring it to update a minimum of one index.

**Budget Analysis** GCG employs a brute force approach, which suffers from significant inefficiencies in both time and space complexity. Specifically, for each iteration, the process includes a gradient calculation that requires a full pass, followed by $B$ substitutions which results in the total number of queries to the model being $T(B+1)$. Each forward pass is quadratic in the overall fixed sequence length $\ell$, and so the overall time complexity is $O(\ell^2 TB)$.

---

**Algorithm 2** Concurrent Greedy Search (CGS)

---

**Input:** Initial prompt $x_{1:n}$, iterations $T$, loss $\mathcal{L}$, threshold $z$, $k$, batch size $B$, suffix indices $J$

$\quad \textbf{repeat } T \textbf{ times}$

$\quad\quad \textbf{for } j \in J \textbf{ do}$

$\quad\quad\quad \Sigma_j = \text{Top-}k(-\nabla_{e_{x_j}}\mathcal{L}(x_{1:n})$ $\hfill \triangleright$ Compute top-$k$ token substitutions

$\quad\quad \textbf{for } b \in \{1, \cdots B\} \textbf{ do}$

$\quad\quad\quad \tilde{x}^{(b)}_{1:n} \leftarrow x_{1:n}$

$\quad\quad\quad I \leftarrow \{\text{Bernoulli}(z)\}_{j \in J}$ $\hfill \triangleright$ Get random indices by taking indices in $J$ with probability $z$.

$\quad\quad\quad \tilde{x}^{(b)}_j := \text{Unif}(\Sigma_j), \forall j \in I$ $\hfill \triangleright$ Perform random substitutions for each index

$\quad\quad x_{1:n} \leftarrow \tilde{x}^{(b^*)}_{1:n}, b^* = \arg\min_b \mathcal{L}(\tilde{x}^{(b)}_{1:n})$ $\hfill \triangleright$ Get best replacement

**Output:** $x_{1:n}$

---

| Dataset | Algorithm | ASR |
|---------|-----------|-----|
| HarryPotterQA | Greedy Coordinate Gradient | 0.87 |
| | Dynamic Suffix Search | 0.17 |
| | Concurrent Greedy Search | 0.90 |
| HarryPotterQA-P | Greedy Coordinate Gradient | 0.83 |
| | Dynamic Suffix Search | 0.2 |
| | Concurrent Greedy Search | 0.87 |

Table 1: Attack Success Rates (ASRs) of different algorithm-dataset pairs.

In contrast, each DSS iteration includes a gradient calculation that requires a forward pass, but only at one index, followed by $B^*$ substitutions for each of $m$ beams. This brings the overall number of queries to the model as $T^*(mB^* + 1)$. However, since $T^*$ in this case is the number of updates made to the model, upper bounded by the sequence length, we have that $T^* << T$ in most cases. While the overall complexity of the model is the same, in practice due to the increasing suffix length and the gradient calculation only in one index, we find that DSS performs attack faster than GCG.

The only differences between GCG and CGS are in which indices to update, so they have the same budget analysis. However, it is expected that CGS converges faster than GCG due to its ability to update multiple indices at the initial optimization steps.

## 4  Results

We run all three of our algorithms on the two datasets HarryPotterQA and HarryPotterQA-P. We determine an attack to be successful if the greedy decoded answer contains the true target answer. We only stipulate that the answer contain the target answer because despite the system instructions, the llama-2-hp model has the tendency to be chatty.

For each algorithm-data pair $(\mathcal{A}, \mathcal{D})$, we report the attack success rate (ASR). Specifically, given a prompt-target pair $(x, y)$, we let $p(x, y, \mathcal{A})$ denote the generated adversarial prompt. We let the indicator function $\mathbb{1} : (p, y) \mapsto \{0, 1\}$ denote if the adversarial prompt $p$ satisfies $d(p, y) > c$ where $c = 0.85$ for HarryPotterQA and $c = 0.95$ for HarryPotterQA-P. In this case, $d$ is similarity function given by the difflib sequencematcher ratio along with the Levenshtein Distance.[2] The reason for the increased threshold value for HarryPotterQA-P is due to the increased sequence length. For lesser threshold values we empirically saw many more false positive matches.

Then we can denote

$$\text{ASR}(\mathcal{A}, \mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(x,y) \sim \mathcal{D}} \mathbb{1}(p(x, y, \mathcal{A}), y)$$

We show the ASRs for each algorithm-dataset pair in Table 1. In addition to ASR, we also give a notion of how efficiently each algorithm is able to retrieve the targeted answers. We greedy decode the updating suffix at various intervals and count how many prompts are able to be successfully attacked with early termination. We omit the results for DSS since all attacks were either successful by first intermediate evaluation step or failed. We show the results in Fig. 1.

## 5  Discussion

### 5.1  Algorithms and Datasets

From the results in Table 1, it is clear that GCG and CGS perform much better than DSS. Our initial assumption was that DSS would be able to leverage the causal nature of Transformer decoders in order to specify the order of which indices get updated to save on computational costs. This assumption takes heavy inspiration from beam search decoding, in which multiple possible generations are decoded in parallel to determine which generation has the highest log-probs.

The crucial difference between DSS and beam search decoding is that in beam search, the possible generations are the exact target of the log-prob maximization whereas in DSS, the possible suffixes
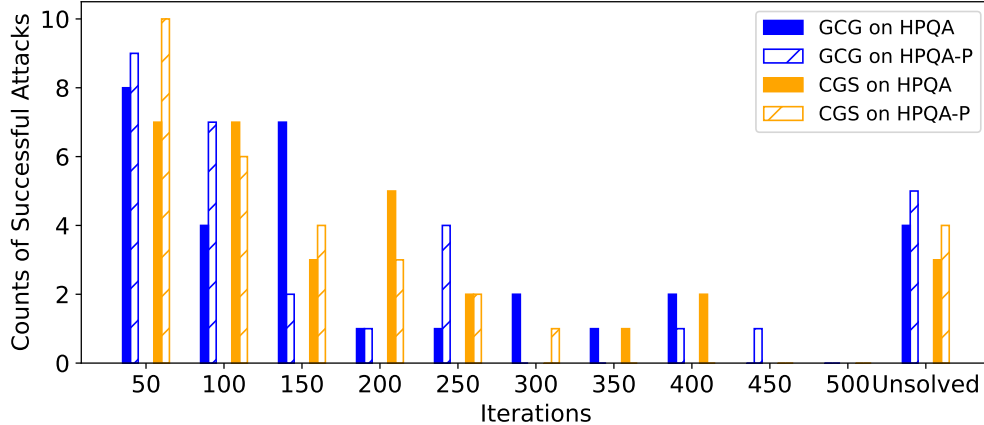
---

[2]https://docs.python.org/3/library/difflib.html

Figure 1: Efficiency of GCG and CGS in attacking `HarryPotterQA` and `HarryPotterQA-P`. Each dataset consists of thirty prompts. If the final answer was incorrect, the prompt is counted as unsolved even if the correct answer was able to be greedily decoded earlier. Otherwise, the instance of the first iteration when the answer is able to be greedily decoded is counted.

serve as more context to the target of log-prob maximization. To make this clear, we see that given a context $x$, beam search and DSS differ in the maximization domain and target.

$$\text{Beam search: } \arg\max_{y_{1:t}} p(y_{1:t} \mid x) \qquad \text{DSS: } \arg\max_{y_{1:t}} p(z \mid [x : y_{1:t}])$$

We can also see this empirically by looking at the greedy decodes of partial suffixes in attacks under DSS. For the prompt "Who was the headmaster of Hogwarts when Harry first arrives," one of the partial suffixes (of length 32) greedy decodes the target answer even though the final suffix does not.

GCG performs worse than CGS in terms of ASR on both `HarryPotterQA` and `HarryPotterQA-P`. The difference is only in one prompt, but interestingly the sets of prompts that were not solved by GCG and CGS differ greatly. We hypothesize that this phenomenon is due to the large amounts of randomization inherent in both algorithms. We explore the differences between these unsolved sets as well as which prompts are the most difficult to attack in Section 5.2.

Besides ASR, we note CGS is able to attack prompts slightly more effectively than GCG. On `HarryPotterQA`, CGS was able to successfully attack 22 prompts within 200 iterations, while GCG was only able to successfully attack 20 prompts. On `HarryPotterQA-P`, CGS was able to successfully attack 23 prompts while GCG was only able to attack 20. However, we note that the randomness in the algorithms likely affects these results as well.

On both GCG and CGS, the ASR on `HarryPotterQA` was slightly higher than on `HarryPotterQA-P`. This is due to the difference in threshold settings for matching decoded answers with true answers. For `HarryPotterQA-P`, a suitable threshold efficiently minimizes false negatives despite lower similarity. In contrast, using the same threshold for `HarryPotterQA` results in more false positives, especially with long sentences differing by one keyword. Therefore, we slightly increased the threshold for `HarryPotterQA`, leading to a higher ASR.

## 5.2 Which prompts are the easiest to adversarially attack?

We introduce the notion of how difficult a prompt is to attack in this section, to see the correlation with difficulty and attack success rate. For a prompt-answer pair $(q, a)$, we define a metric $m : (q, a) \mapsto r \in [0, 1]$ where

$$r = \frac{\log p_{\texttt{llama2-chat}}(a \mid q)}{\log p_{\texttt{llama2-hp}}(a \mid q)}$$

The closer $r$ is to 1 indicates that the log-probs of the answer was not changed by the fine-tuning of `llama2-chat` to `llama2-hp`. Lower values of $r$ indicate that the log-probs of the answer were drastically changed by fine-tuning, so we expect the difficulty in attacking these prompts to be greater.
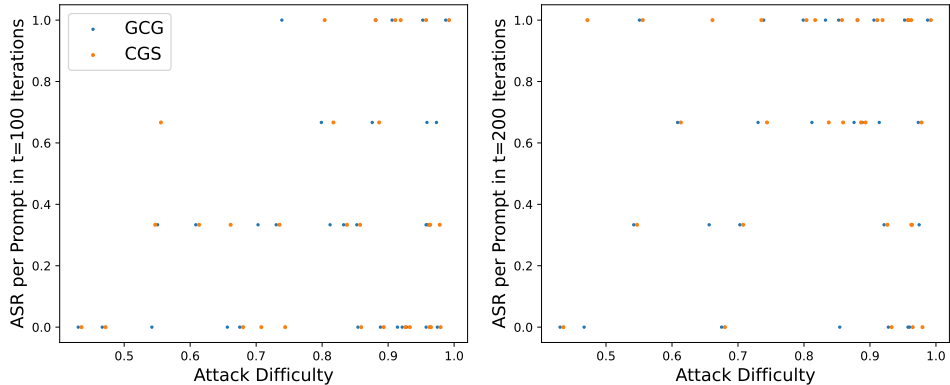
Figure 2: The attack success rates of prompts plotted against their difficulties. Success rates per prompt were obtained by running the algorithm on the same prompt-answer pair three times. Attack difficulties were calculated using the metric described above.

As a motivation to this metric, we note that the model `llama2-hp` unlearned the Harry Potter content through decreasing the log-probs of tokens that are classified as being Harry Potter content, such as character names or important items to the plot. Because the unlearning techniques in the fine-tuning were quite shallow (the fine-tuning only ensured that the weights in the final layer do not produce tokens related to Harry Potter), we believe that the change in the log-probs of the unlearned model from the base `llama2-chat` model will act as an indicator of how much this knowledge was erased from the model.

For our experiments, we run GCG and CGS on `HarryPotterQA` for 100 and 200 iterations, repeating the attacks three times. For each prompt, we plot the ASR for that prompt (the number of times that prompt was successfully attacked divided by three) and the difficulty. We show the results in Fig. 2.

Despite there being some outliers with high $r$ values that are were never successfully attacked, and some with low $r$ values that were consistently attacked (but only under CGS), we find that the general trend is that prompts that deviate more from their base log-probs are more difficult to attack.

Moreover, these results give us a better means of comparing GCG and CGS. We find that over the 90 prompts attacked (3 copies of the original 30), CGS is able to successfully attack at least 10% more prompts than GCG in 200 iterations (58 vs 52).

| Prompt | Greedy Decoded Output |
|---|---|
| Who are Harry Potter's best friends? | Sarah, Emily, and Jack. |
| *+ presumably XIV #( wyseauescape XII marker !germ p\ u0142 \%, ernal Lisa.... fonce tell przew separator* | Ron Weasley and Hermione Granger. |
| What house is Harry Potter sorted into at Hogwarts? | Thornwood. |
| *+ gentlemenG!– Lou comport fif from Australeff gepubliceerd']; cour fris D sorts Chfif inde* | Gryffindor. |
| What object does Harry catch to end a Quidditch match? | The ball. |
| *+ emphas oro material( Pol earliest Ŏ0e9t ele theyyyypg dil Goldenpel Only this sorte XII worthy Affairsurrey* | The Golden Snitch. |

Table 2: Greedily decoded outputs from `llama-2-hp` with the given prompt inputs. The adversarial suffixes are concatenated with the prompt to produce the correct answers under greedy decoding.
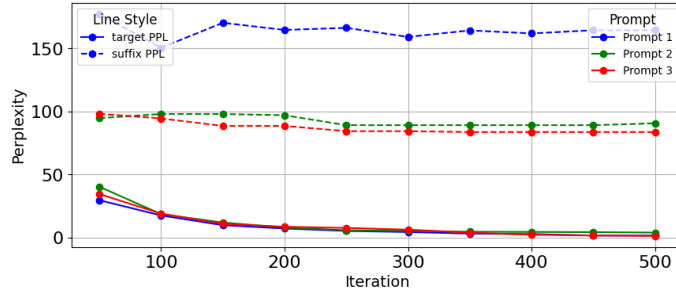
Figure 3: The perplexity trends of the targets and suffixes throughout the attack for three prompts. Dashed lines indicate the perplexities of the tokens in the suffix, solid lines indicate the perplexities of tokens in the target. The different colors indicate different prompts.

## 5.3 How interpretable are the adversarial suffixes?

To assess the interpretability of adversarial suffixes, we selected three examples for examination, as presented in the Table 2. Due to their length, these suffixes are truncated in the table, but it remains evident that they are generally unreadable. The character strings do not form coherent words or phrases that could be understood or used to answer questions based on common human knowledge. Nonetheless, they are successful in inducing the unlearned model in generating supposedly forgotten information.

To better understand the nature of adversarial suffixes, we analyzed the perplexity values associated with these suffixes over a series of iterations, comparing them against the target suffix. Our experiment shows a consistent decrease in the perplexity of the target suffix as training progresses. In contrast, the perplexity of the adversarial suffixes remains high throughout the attack phase. Thus we conclude it is hard to interpret adversarial suffixes generated by the greedy algorithm. The result is in Fig. 3. These results were obtained by running GCG on the three prompts on the `HarryPotterQA` dataset.

## 6 Conclusion

Our project focuses on the auditing of unlearned language models through adversarial prompt generation. We show that the model that unlearned Harry Potter content, referred to as `llama2-hp`, still retains supposedly forgotten information (Eldan & Russinovich, 2023). We compare three algorithms in our experiments: GCG, introduced by Zou et al. (2023) for adversarially prompting in the domain of breaking RLHF filters; DSS, our algorithm and a novel approach towards suffix generation that seeks to leverage the causal nature of decoder LLMs to perform fast attacks which are linear in the suffix length; CGS, our algorithm, a direct modification of GCG which allows for quicker convergence due allowing multiple indices to be updated in every iteration. We ultimately find that our algorithm CGS outperforms GCG by over 10%, with the same computational cost of attack. DSS performed its attacks much quicker but its effectiveness trailed GCG and CGS immensely.

We also perform further analyses on the interpretability of the generated suffixes as well as introducing a metric quantifying how hard prompts are to attack.

## 7 Contributions and Reproducibility

Our code can be found at this repository. Refer to the README on the master branch for examples on how to run the attacks. Requirements are also listed in the README.

`Jeffrey`    I wrote the base attack classes and implemented the the attacks in attacks.py, as well as the code to parse config files for attack parameters in run_attack.py. In the report, I contributed to the Introduction, Related Work, the Adversarial Prompting subsection, and results sections.

`Steven`    I generated questions in `GPT-4` and prompted it into `llama-2-chat` to form our dataset, and also classify correct and wrong qa pairs for the dataset. In the report, I contributed to Dataset and Time & Space complexity analysis.

# References

Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani Srivastava, and Kai-Wei Chang. Generating natural language adversarial examples, 2018.

Battista Biggio, Igino Corona, Davide Maiorca, Blaine Nelson, Nedim Šrndić, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. *Evasion Attacks against Machine Learning at Test Time*, pp. 387–402. Springer Berlin Heidelberg, 2013. ISBN 9783642387098. doi: 10.1007/978-3-642-40994-3_25. URL http://dx.doi.org/10.1007/978-3-642-40994-3_25.

Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks, 2017.

Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models, 2023.

Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. Hotflip: White-box adversarial examples for text classification, 2018.

Ronen Eldan and Mark Russinovich. Who's harry potter? approximate unlearning in llms, 2023.

William Fleshman, Aleem Khan, Marc Marone, and Benjamin Van Durme. Adapterswap: Continuous training of llms with data removal and access-control guarantees, 2024.

Antonio Ginart, Melody Y. Guan, Gregory Valiant, and James Zou. Making ai forget you: Data deletion in machine learning, 2019.

Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples, 2015.

Jonathan Hayase, Ema Borevkovic, Nicholas Carlini, Florian Tramèr, and Milad Nasr. Query-based adversarial prompt generation, 2024.

Robin Jia and Percy Liang. Adversarial examples for evaluating reading comprehension systems, 2017.

Erik Jones, Anca Dragan, Aditi Raghunathan, and Jacob Steinhardt. Automatically auditing large language models via discrete optimization, 2023.

Gaoyang Liu, Xiaoqiang Ma, Yang Yang, Chen Wang, and Jiangchuan Liu. Federated unlearning, 2021.

Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks, 2016.

Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. Scalable extraction of training data from (production) language models, 2023.

Ayush Sekhari, Jayadev Acharya, Gautam Kamath, and Ananda Theertha Suresh. Remember what you want to forget: Algorithms for machine unlearning, 2021.

Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models, 2024.

Taylor Shin, Yasaman Razeghi, Robert L. Logan IV au2, Eric Wallace, and Sameer Singh. Autoprompt: Eliciting knowledge from language models with automatically generated prompts, 2020.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti 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, 2023.

Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. Universal adversarial triggers for attacking and analyzing nlp, 2021.

Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, Sang T. Truong, Simran Arora, Mantas Mazeika, Dan Hendrycks, Zinan Lin, Yu Cheng, Sanmi Koyejo, Dawn Song, and Bo Li. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models, 2024.

Heng Xu, Tianqing Zhu, Lefeng Zhang, Wanlei Zhou, and Philip S. Yu. Machine unlearning: A survey, 2023.

Jingwen Ye, Yifang Fu, Jie Song, Xingyi Yang, Songhua Liu, Xin Jin, Mingli Song, and Xinchao Wang. Learning with recoverable forgetting, 2022.

Dawen Zhang, Pamela Finckenberg-Broman, Thong Hoang, Shidong Pan, Zhenchang Xing, Mark Staples, and Xiwei Xu. Right to be forgotten in the era of large language models: Implications, challenges, and solutions, 2023.

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