



ALPACA AGAINST VICUNA: Using LLMs to Uncover Memorization of LLMs

Aly M. Kassem^{1*} Omar Mahmoud^{2*} Niloofar Mireshghallah^{3*}
Hyunwoo Kim⁴ Yulia Tsvetkov³ Yejin Choi⁵ Sherif Saad¹ Santu Rana²

¹University of Windsor ²A2I2, Deakin University

³University of Washington ⁴NVIDIA ⁵Stanford University

kassem6@uwindsor.ca, o.mahmoud@deakin.edu.au, niloofar@cs.washington.edu

Abstract

In this paper, we investigate the overlooked impact of instruction-tuning on memorization in large language models (LLMs), which has largely been studied in base, pre-trained models. We propose a black-box prompt optimization method where an *attacker* LLM agent uncovers higher levels of memorization in a *victim* agent, surpassing traditional approaches that prompt the model directly with training data. Using an iterative rejection-sampling process, we design instruction-based prompts that minimize overlap with training data to avoid providing direct solutions while maximizing overlap between the victim’s output and the training data to induce memorization. Our method shows 23.7% more overlap with training data compared to state-of-the-art baselines. We explore two attack settings: an analytical approach that determines the empirical upper bound of the attack, both with and without access to responses for prompt initialization, and a practical classifier-based method for assessing memorization without access to memorized data. Our findings reveal that instruction-tuned models can expose pre-training data as much as, or more than, base models; contexts beyond the original training data can lead to leakage; and instructions generated by other LLMs open new avenues for automated attacks, which we believe require further exploration.¹


1 Introduction

Pre-trained language models are commonly instruction-tuned for user-facing applications to generate high-quality responses to task-oriented prompts (Ouyang et al., 2022; Taori et al., 2023a; Chowdhery et al., 2023). While extensive prior work has investigated memorization in pre-trained base LLMs and its implications for privacy, copyright, and fairness (Carlini et al., 2022; Bider-

man et al., 2023a; Shi et al., 2023; Mireshghallah et al., 2022), there is limited understanding of how instruction-tuning affects the memorization and discoverability of pre-training data in aligned models. Studies have shown that aligned LLMs can emit training data up to 150× more often than in regular operation (Nasr et al., 2023). To address this gap, we pose the question: *Can we use instruction-based prompts to uncover higher levels of memorization in aligned models?* The established method of quantifying memorization (Carlini et al., 2023) assumes that a sequence d is memorized if prompting the model with the original prefix from the training data yields sequence d (or a similar sequence for approximate memorization; Biderman et al. 2023a). However, recent findings suggest that prompts other than the original training data may trigger even higher levels of regurgitation (Schwarzschild et al., 2024). To explore this, we propose a new optimization method, illustrated in Figure 1, where an aligned language model acts as an ‘attacker,’ generating prompts that induce a victim (target) model to produce outputs more faithful to the training data. The attacker refines prompts through a feedback loop guided by a reward function that increases the overlap between the victim’s output and the ground truth. This approach is inspired by adversarial methods in computer security literature (Wang et al., 2023a) and has been effective in jailbreaking attacks (Mehrotra et al., 2023a; Zeng et al., 2024; Ramesh et al., 2024).

To evaluate our approach, we draw parallels between safety jailbreaking techniques and training data extraction, using automatic prompt optimization to guide models toward outputs aligned with their training data. However, unlike jailbreaking, our goal is not to bypass specific safety features but to examine memorization. We evaluate our method using Greedy Coordinate Gradient (CGC; Zou et al. 2023), a white-box prompt optimization technique,

*Equal Contribution

¹ https://github.com/Alymostafa/Instruction_based_attack

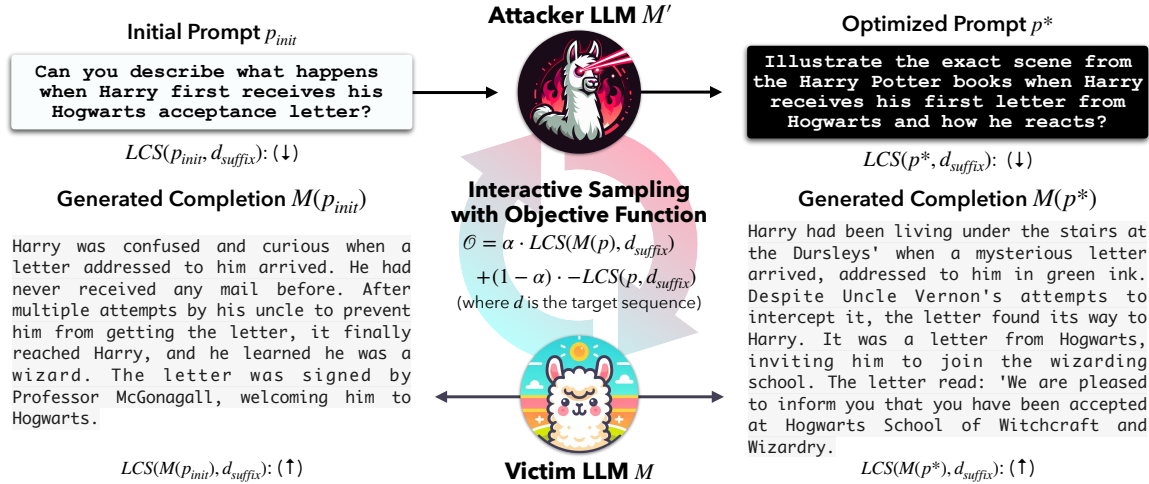


Figure 1: Overview of our method: we first create an initial prompt that turns the target training sequence into an instruction. The attacker LLM uses this prompt to generate multiple candidate prompts designed to make the victim LLM produce responses that closely match the training data. We score each candidate based on two criteria: (1) the overlap between the victim’s response and the training data (higher is better) and (2) the overlap between the candidate prompt and the training data (lower is better to avoid revealing the solution). This score guides the attacker in optimizing and generating new prompts for further rounds of optimization.

and compare it to methods like Reverse-LM (Pfau et al., 2023) and sequence extraction (prefix-suffix; Carlini et al. 2022, 2021) across both base and instruction-tuned models. Our method was tested on Llama-based, OLMo, and Falcon models (Touvron et al., 2023; Penedo et al., 2023; Groeneveld et al., 2024), and their instruction-tuned variations, such as Alpaca (Taori et al., 2023a), Tulu (Wang et al., 2023b), and Vicuna (Chiang et al., 2023b), using sequences of 200, 300, and 500 tokens from five pre-training data domains (Duan et al., 2024). We find that our approach uncovers 23.7% more memorization in instruction-tuned models compared to the prefix-suffix method (Carlini et al., 2022), which can give a false sense of privacy. Furthermore, our method reveals 12.4% higher memorization in instruction-tuned models, indicating that contexts beyond the original pre-training data can lead to leakage, highlighting the need for improved privacy measures. To demonstrate the real-world applicability of our method, we conduct four case studies: regurgitation of copyrighted material subsection 6.1, privacy auditing of LLMs subsection 6.2, refusal behavior of LLMs subsection 6.3, and the development of a classifier that detects whether a prompt can elicit memorized data without needing access to the response or the memorized content, enabling a more practical attack subsection 6.4. Our method achieved 39% more extraction in copyright-related queries on the Books3, BooksMIA, and NYT datasets (Com-

puter, 2023; Shi et al., 2023; Grynbaum and Mac, 2023), and a 56.6% increase in privacy auditing over the prefix-suffix approach (Eldan and Ruvinsky, 2023). Additionally, we show that LLMs do not refuse copyright-related queries with our approach, demonstrating high adversarial effectiveness. Lastly, our classifier reliably detects prompts triggering memorized data in our framework without requiring the actual response, proving more practical. We hope these results encourage further research into automated model auditing and probing using LLMs to develop more efficient data reconstruction methods.

2 Background: Quantifying Memorization

In this work, we use the discoverable notion of memorization for LLMs and quantify it through approximate string matching. Below, we define these terms.

Definition 1 (Discoverable Memorization) An example $x = [p||s]$, drawn from training data D , is considered memorized by model f_θ if $f_\theta(p) = s$, where x consists of a prefix p and a corresponding suffix s .

The concept entails that the prefix guides the model’s generation process towards the most probable completion, typically the suffix if the example has been memorized. Drawing from previous research, Carlini et al. (2022) identified certain fac-

tors significantly influencing memorization, including model size, utilization of data deduplication techniques, and contextual aspects.

Definition 2 (Approximate String Matching)

For a model f_θ and a given similarity metric β , an example x from the training data D is said to be approximately memorized if there exists a prompt p such that the output of the model $f_\theta(p)$ is s' , where s and s' are close in accordance with the similarity metric β , i.e., $\beta(s, s')$ is high.

Prior research demonstrates approximate memorization’s superiority over verbatim memorization in LLMs (Ippolito et al., 2023; Biderman et al., 2023a). We employ ROUGE-L to measure the similarity via the longest common subsequence between model-generated and original continuations, adhering to approximate memorization in our work.

3 Using LLMs to Probe Memorization in other LLMs

In this section, we begin by formally outlining the optimization problem and specifying our objective function. We present our method’s pipeline, as shown in Figure 1 and Algorithm 1, which includes initialization, sampling, and refinement, creating the optimized prompt.

3.1 Problem Formulation

Consider a set of sequences $D = \{d_1, \dots, d_N\}$, where D is the pre-training dataset of the LLM model M . A function $f : d \rightarrow p^*$ is a transformation process that takes a pre-training sequence $d \in D$ and generates an optimized prompt p^* that maximizes the overlap between the output sequence of the model $M(p^*)$ and pre-training sequence d :

$$p^* = \arg \max_p \mathcal{O}_{d,M}(p)$$

where $\mathcal{O}_{d,M}(p) = \text{LCS}(M(p), d_{\text{suffix}})$ is the objective function to maximize for a fixed model M and sequence d . $M(\cdot)$ denotes the operation of decoding from the model M , conditioned on a given input. LCS is the longest common subsequence that measures the syntactic similarity between sequences, and in our case, we employ ROUGE-L (Lin, 2004).

Algorithm 1 Interactive Sampling Algorithm

```

1: Input: pre-training sample  $d, M, M', M_{\text{init}}$ 
2:  $p_{\text{init}} \leftarrow M_{\text{init}}(d)$  //Construct initial prompt
3:  $p_{t-1} \leftarrow p_{\text{init}}$ 
4: for  $t = 3$  do
5:    $p_t \sim M'(Instr|p_{t-1}, n = 24)$  //Sample 24
6:    $\mathcal{O} = \alpha \cdot \text{LCS}(M(p_t), d_{\text{suffix}}) + (1 - \alpha) \cdot$ 
      $-\text{LCS}(p_t, d_{\text{suffix}})$ 
7:    $p_t = \arg \max(\mathcal{O})$  //Obtain the highest scoring prompt
8: end for
9:  $p^* = \arg \max(p_0, \dots, p_t)$  //get the highest over iters
10: return  $p^*$  //Return optimal prompt

```

We consider two settings for the proposed attack. The first one is to estimate its empirical upper bound, which is the default assumption throughout the paper. The second one is a practical setting where we don’t use the full sequence either for evaluation or initialization.

1) Empirical Upper-Bound. To better estimate the empirical upper bound of the attack, we assume that we have access to the full sequence d , where sequence d is split into d_{prefix} and d_{suffix} . We use the full sequence for initialization, which will be discussed later, and for feedback in the objective function, which can be directly used to maximize $\text{LCS}(M(p), d_{\text{suffix}})$. However, LLMs have been shown to regurgitate and repeat their inputs (Zhang and Ippolito, 2023; Priyanshu et al., 2023). Therefore, an obvious solution could be $p = [z||d]$, where z is an instruction like "repeat". To avoid this shortcut, we rewrite the objective \mathcal{O} as follows to de-incentivize such solutions:

$$\mathcal{O} = \alpha \cdot \text{LCS}(M(p), d_{\text{suffix}}) + (1 - \alpha) \cdot (-\text{LCS}(p, d_{\text{suffix}}))$$

We include the second term to penalize solutions significantly overlapping with the sequence d_{suffix} . The hyperparameter α regulates how much d is utilized, balancing a high memorization score with minimal overlap with the ground truth (see Appendix A for details).

2) Practical Setting.

In practical scenarios where the suffix d_{suffix} is inaccessible thus, we can not utilize ROUGE-L for feedback. As a result, we use the d_{prefix} only for prompt initialization and evaluation. We learn a function $C : \mathcal{P} \rightarrow \mathcal{L}$, where C is a binary classifier that takes a prompt $p \in \mathcal{P}$ and outputs a label $l \in \mathcal{L}$, where $\mathcal{L} = \{T, NT\}$ represents the possibility that a prompt would trigger memorized responses or not. Assume we have access to preference data

$\mathcal{D}_{\text{pref}} = \{(p, l) \mid p \in \mathcal{P}, l \in \mathcal{L}\}$. We will discuss the details of the classifier in [subsection 6.4](#).

3.2 Optimization via Interactive Sampling

Initialization. To create the initial prompt, the training data point is transformed into a question. We consider two setups where we use the full sequence d or the prefix only d_{prefix} (see [section 7](#)). An initialization function

$$I : \{d_{\text{prefix}}, (d_{\text{prefix}}, d_{\text{suffix}})\} \rightarrow P_{\text{init}}$$

is defined, where $\{d_{\text{prefix}}, (d_{\text{prefix}}, d_{\text{suffix}})\}$ represents either the prefix alone or both the prefix and suffix.

We instruct LLM with a ‘meta-prompt’ along with the pre-training sample. We also add customized instructions to regularize the prompts to keep them abstract and not overly lengthy. We use the meta-prompt on GPT-4 to help generate the initial prompt. Still, we show that utilizing other models, such as Mixtral ([Jiang et al., 2024](#)), also yields comparable performance ([section 7](#)).

Interactive Loop. Upon receiving the initial prompt, we employ a two-step strategy to optimize it for the best results, involving *exploration* and *exploitation*.

In our setting, we use an alternate model $M'(\cdot, |\text{instr}|)$, with a specific instruction instr , as an attacker model that proposes prompts p . We perform constrained sampling $p_t \sim M'(\cdot, |\text{instr}||p_{t-1}|)$ at time step t from the proposal distribution, where the constraint is to maximize $\text{LCS}(M(p_t), d_{\text{suffix}})$. This is achieved with rejection sampling (best-of- n) from M' .

(1) **Best-of- n sampling from M' :** During optimization, the meta-prompt text evolves from its initialization. We instruct the model to paraphrase the previous prompt p_{t-1} and generate a new one. The attacker LLM produces 24 new prompts per sample, which are scored using our objective function (ROUGE if suffix access is available, otherwise via the proposed classifier) as shown in steps 5, 6, and 7 in [Algorithm 1](#). The highest-scoring prompt is selected, ensuring better-quality samples in the next step where we employ refinement.

(2) **Refine:** To proceed, We designate the improved prompt from the previous iteration as the starting point and repeat the sampling process three times, following step 4 in [Algorithm 1](#). Each iteration incorporates feedback from the victim to refine the prompt, thereby enhancing ex-

traction capabilities and engaging with the attacker LLM using the previous prompt. At time step t , we apply constrained sampling $p_t \sim M'(\cdot, |\text{instr}, ||, p_{t-1}|)$, where the constraint is maximizing $\text{LCS}(M(p_t), d)$, using rejection sampling (best-of- n) from M' .

4 Experimental Settings

4.1 Attacker & Victim LLMs

Attacker LLMs: Our method leverages the open-source Zephyr 7B model, an instruction-tuned variant of Mistral-7B β ([Tunstall et al., 2023](#)), as the attacker due to its exceptional ability to follow instructions and generate text effectively at the time of writing this paper. We also showcase employing more powerful LLMs as attackers (e.g. GPT-4) in [section 7](#).

Victim LLMs: We assess the memorization capabilities of instruction-tuned LLMs compared to their base model across various sizes (7B, 13B, 30B) by applying our method on five open-source models of different sizes by employing the instruction-tuned versions of Llama-1 (Alpaca, Tulu, Vicuna) ([Touvron et al., 2023](#); [Taori et al., 2023b](#); [Wang et al., 2023b](#); [Chiang et al., 2023a](#)), OLMo ([Groeneveld et al., 2024](#)), and Falcon ([Penedo et al., 2023](#)) since there is a disclosure in their training data. By comparing these instruction-tuned models to their base model, we gain insights into the impact of instruction-tuning on memorization. See [Appendix D](#) for more details about the models.

4.2 Evaluation Data

Data Domains: We construct diverse evaluation datasets by sampling from several pre-training datasets used in base models. Specifically, we use Llama (replicated from RedPajama due to data unavailability), Falcon’s RefinedWeb (from Common Crawl), and OLMo’s Dolma. Llama spans five domains (C4, CC, Arxiv, Books, and Github), while Dolma covers six domains (C4, CC, Arxiv, Books, Reddit, Stack, and PeS2o). We ensure uniformity in sequence length distribution, selecting 15,000 samples from Llama, 3,000 from Falcon’s RefinedWeb, and 16,000 from OLMo’s domains.

Sequence Lengths Selection: To evaluate adaptability across different sequence lengths (200, 300, and 500), abbreviated as "seq.," we adopt a splitting ratio inspired by real-world usage patterns. Based

Average Over Three Sequence Lengths (200, 300, 500)																
Model	Method	Github			ArXiv			CC			C4			Books		
		Mem ↑	LCS _P ↓	Dis ↑	Mem ↑	LCS _P ↓	Dis ↑	Mem ↑	LCS _P ↓	Dis ↑	Mem ↑	LCS _P ↓	Dis ↑	Mem ↑	LCS _P ↓	Dis ↑
Alpaca	P-S-Inst	.270	.124	-	.179	.112	-	.155	.104	-	.143	.114	-	.131	.093	-
	Reverse-LM	.229	.200	.864	.133	.196	.848	.113	.186	.843	.110	.181	.834	.122	.142	.865
	Ours	.322	.102	.864	.228	.108	.848	.214	.096	.830	.203	.090	.834	.221	.079	.865
Vicuna	P-S-Inst	.273	.125	-	.213	.112	-	.205	.114	-	.191	.114	-	.198	.093	-
	Reverse-LM	.255	.200	.864	.200	.196	.848	.173	.186	.830	.173	.181	.834	.166	.142	.865
	Ours	.325	.096	.864	.232	.104	.853	.213	.092	.838	.201	.084	.841	.223	.079	.866
Tulu	P-S-Inst	.274	.124	-	.207	.112	-	.170	.106	-	.137	.114	-	.172	.093	-
	Reverse-LM	.245	.200	.864	.153	.196	.848	.121	.186	.830	.117	.181	.834	.135	.142	.865
	Ours	.359	.104	.857	.237	.104	.851	.221	.094	.835	.210	.086	.836	.233	.079	.865
Seq Len	Tulu-7B															
200	P-S-Inst	.298	.125	-	.216	.107	-	.176	.103	-	.140	.111	-	.188	.090	-
	Reverse-LM	.254	.191	.877	.154	.200	.890	.130	.203	.863	.123	.195	.862	.153	.151	.880
	Ours	.372	.098	.877	.204	.093	.883	.225	.104	.858	.214	.095	.853	.236	.082	.882
300	P-S-Inst	.276	.124	-	.209	.112	-	.174	.106	-	.142	.114	-	.178	.095	-
	Reverse-LM	.246	.203	.881	.157	.196	.853	.125	.190	.822	.116	.182	.826	.134	.145	.877
	Ours	.341	.084	.878	.248	.108	.856	.222	.099	.824	.209	.090	.825	.231	.079	.872
500	P-S-Inst	.247	.124	-	.195	.117	-	.159	.102	-	.128	.117	-	.149	.095	-
	Reverse-LM	.233	.204	.833	.147	.192	.803	.107	.164	.805	.112	.167	.814	.118	.129	.838
	Ours	.363	.129	.814	.260	.112	.809	.216	0.079	.824	.207	.074	.829	.231	0.076	.841

Table 1: Comparison of our method with baselines across pre-training data domains. Mem denotes the memorization score (ROUGE-L), LCS_P is input prompt and suffix overlap, and Dis is optimized vs. initial prompt distance. Results are averaged over three sequence lengths on top, and for the *Tulu-7B* model, we show a breakdown at the bottom. The highest performance within each domain is bolded.

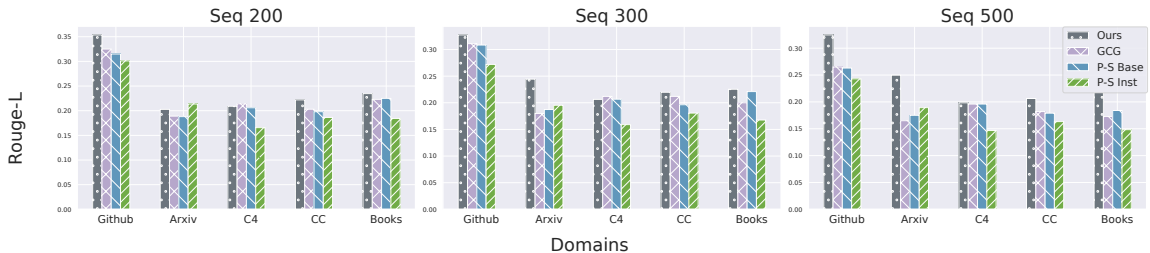


Figure 2: Comparison of our method to the GCG, P-S baseline, and P-S-instruction on the Llama and its instruction-tuned versions. We evaluate different subsets of the pre-training data and observe that our method consistently outperforms the GCG and prefix-suffix baseline.

on analysis from the WildChat dataset (Zhao et al., 2024), we divide each sample, allocating 33% as the prefix and 67% as the suffix, reflecting typical usage scenarios (see Appendix D for further details).

4.3 Baseline Methods

We compare against three methods under two access settings: white box and black box.

(1) **Prefix-Suffix (P-S) sequence extraction (Carlini et al., 2022, 2021)**: A black-box attack where the model is prompted with the first n tokens (prefix) of a pre-training sample to generate output, applied to both base and instruction-tuned models.

(2) **GCG (Zou et al., 2023)**: A white-box adversarial attack that starts with the original prefix and

is trained for thirty epochs on the base model.

(3) **Reverse LM (Pfau et al., 2023)**: A method that reverses token order during training, using a Pythia-160M model trained on the deduplicated Pile dataset (Pfau et al., 2023; Biderman et al., 2023b; Gao et al., 2020).

4.4 Evaluation Metrics

Measuring Memorization/Reconstruction: In our evaluation, we use ROUGE-L to measure memorization by comparing the longest common subsequence between the generated and original suffixes, closely aligning with the memorization score introduced by Biderman et al. (2023a), which emphasizes ordered token matches between model-generated continuations and the true text. **To evalu-**

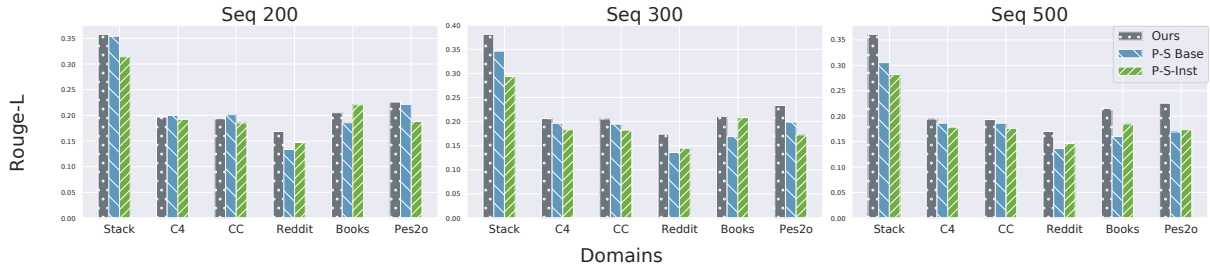


Figure 3: Comparison of our method to the P-S baseline on the OLMo model. We evaluate different subsets of the pre-training data, Dolma, and observe that our method outperforms the prefix-suffix baseline consistently.

ate prompt overlap, particularly in our analytical solution where the prompt includes the ground truth suffix, we assess the overlap between the prompt and suffix to ensure it does not exceed the overlap in the original prefix-suffix combination. We denote this overlap as LCS_p and use ROUGE-L to quantify it.

5 Main Results

Evaluating on Instruction-Tuned LLMs. Table 1 summarizes our main findings and compares them with baselines across different pre-training data domains. Our method reveals significantly higher levels of memorization compared to traditional prefix-suffix methods. On average, our approach achieves a 5% increase in memorization, reaching up to 12% in scenarios with a sequence length of 500. For instance, GitHub & Tulu LM achieve a reconstruction Rouge-L score of 24.7% with prefix-suffix, whereas our method improves this to 35.9%. These results hold consistently across various models, including Llama-based models, OLMo (Groeneveld et al., 2024), and Falcon (Penedo et al., 2023), as well as larger models like 13B and 30B. Detailed results on the Falcon model and larger sizes are provided in Appendix B.

Evaluating on Base LLMs. Figure 2 compares Base and Instruction-tuned LLMs, GCG, and our method. Comparing P-S-Inst and P-S-Base alone would *misleadingly suggest that instruction-tuned models uncover less training data*. However, our method uncovers more memorization than all other baselines, including the base model, showing that instruction-tuned models can reveal more pre-training data when prompted correctly. While the white-box GCG uncovers 1% more memorization than P-S attacks, it still falls short of our method. ReverseLM performs the worst due to its transferability setting from the Pythia model. For detailed results and improvement percentages, refer

to Appendix B. Hyperparameter details are in Appendix A, and optimized prompts and outputs are in Appendix B. For runtime details of the proposed method and GCG, see Appendix D.

Prompt Overlap Analysis. As shown in Table 1, consistently, our method achieves equivalent or lower overlap (LCS_p) in terms of ROUGE-L, with the prefix-suffix baseline. For example, our approach has significantly lower overlap in domains like GitHub, ensuring a fair comparison with baseline methods.

6 Alpaca Vs Vicuna In The Wild

6.1 CASE STUDY: Extraction of Copyrighted Books/Articles

We applied our prompt optimization technique to extract copyright infringements in training data, targeting excerpts from copyrighted books and articles across various models.

Evaluation Data. We used the Books3 subset from the Redpajamas dataset to assess Llama instruction-tuned LLMs and Project Gutenberg to evaluate OLMo, as detailed in subsection 4.2. Additionally, we selected 200 samples from BookMIA and 100 from New York Times articles to evaluate GPT-4o, which has previously been shown to memorize data from these models (Shi et al., 2023; Grynbaum and Mac, 2023).

Results. Figure 3 and Figure 4 demonstrates that our method consistently outperforms Prefix-Suffix in OLMo & Llama based models in Book domain. For GPT-4o, Although it often refuses or avoids verbatim repetition of training data, we achieved approximately 25% overlap on average—doubling the result of simply asking or continuing the text in BookMIA & NYT.

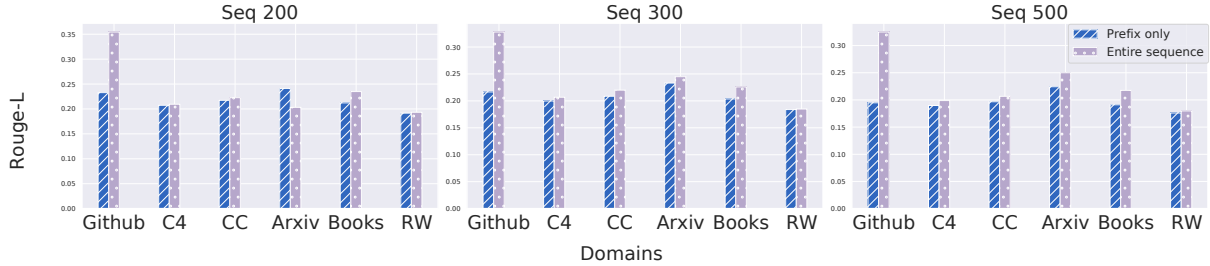


Figure 4: Comparison of our attack performance shows that optimizing prompts over partial sequence access versus full access (default assumption through the paper) shows similar results across domains. This highlights the robustness of optimizing prompts with limited sequence information.

6.2 CASE STUDY: Eliciting Unlearned Harry Potter

(Eldan and Russinovich, 2023) introduced an unlearning technique to remove knowledge of the Harry Potter books through multiple unlearning steps on Llama-2 (Touvron et al., 2023), resulting in a model that no longer retains the targeted content. Although querying the model before and after unlearning shows it has forgotten the information, we aim to assess the model’s behavior under *adversarial prompts* using our approach.

Evaluation Data. We sampled 300 passages from various Harry Potter books, each with a sequence length of 300 to provide sufficient context for prompt generation.

Results. Our optimized prompts elicited highly similar completions to the original text, achieving 23.6% overlap compared to 10.2% using prefix-suffix prompts. These findings suggest the unlearning technique is vulnerable to adversarial prompts that deviate from the original training context.

6.3 CASE STUDY: LLMs Refusal

OpenAI models frequently refuse to answer certain questions, particularly those that seek harmful responses, such as inquiries about illegal activities or hate, harassment, and violence. Recently, when prompted to continue a passage from a book or article, these models declined to respond.

Evaluation Data. We use the same Harry Potter book subset from unlearning to assess refusal rates, with GPT-4o as the judge. GPT-4 and GPT-4o were evaluated on the prefix-suffix and our generated prompts. We also assessed overlap and ensured the completions closely matched the ground truth from our prompts.

Results. By comparing our generated prompts with the prefix-suffix method, we found that our approach bypasses filters, yielding responses for all 300 samples, while the prefix-suffix refusal rates

are 13.65% for GPT-4 and 26.19% for GPT-4o. This demonstrates the robustness of our method in adversarial generation.

6.4 CASE STUDY: Predicting Memorization For Practical Attack

We previously discussed the applicability of our approach when d_{suffix} is inaccessible, as often occurs in real-world scenarios. We developed a classifier to replace the ROUGE-L function in our optimization loop to address this. We outline the preference data and then investigate the classifier’s technical details.

Preference Data. We ran several iterations using the full training sequence, collecting optimized and non-optimized prompts per sample. Each iteration produced one optimized and 23 non-optimized prompts, generating 24 samples over three iterations. We unified preference data by merging sequence lengths and victims to train a single classifier per domain, and we downsampled non-optimized classes to overcome the data imbalance problem.

Technical Details. We train a single classifier for each domain, encompassing various sequence lengths and target entities, using DeBERTa-v3-large (He et al., 2021) with weighted CrossEntropy loss. The model is trained on an H100 80GB GPU for 1500 steps with a batch size of 16 and a maximum sequence length of 512. The dataset of 20,000 samples is split into 80% training, 10% validation, and 10% testing.

Results. We assess the classifier’s performance using the macro F1 score across different data domains, achieving an average F1 score of 70% in distinguishing prompts that trigger memorized responses from those that do not. While the classifier’s performance may not be optimal, we consider this a significant step toward practical attacks in future work, which could be improved by integrating

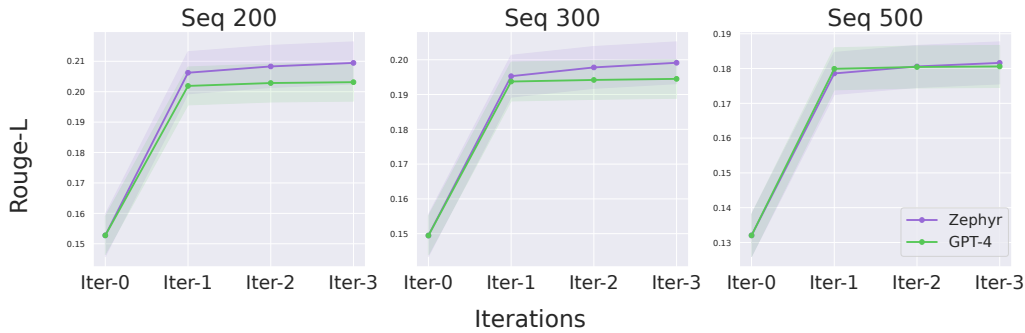


Figure 5: A comparison of our method’s performance using Zephyr and GPT-4 as attacker LLMs is shown for different iteration steps during optimization. We observe that the performance increases across varying sequence lengths as optimization iterations increase.

prefixes with questions in the Natural Language Inference (NLI) task or utilizing Direct Preference Optimization (DPO) (Rafailov et al., 2024) to better align an LLM with the distribution of optimized prompts.

7 Ablation & Analysis

In this section, we conduct ablations and analyses to identify the key components that drive our method’s improvements over baseline models.

GPT-4 is NOT the best attacker. We evaluate GPT-4 as an alternative attacker and find that Zephyr consistently outperforms GPT-4 at a sequence length of 200, maintaining a margin of 0.05 across all domains, as shown in Figure 5. While the performance gap narrows at a sequence length of 300, Zephyr still leads. At 500 tokens, however, GPT-4 begins to match or exceed Zephyr, particularly in the ArXiv domain, where summarization complexity increases with longer sequences.

Initialization without Suffix. In previous experiments, we used the full training sequence, including suffixes, to test Instruction-Tuned LLMs with an overlap penalty to prevent cheating. In real-world settings, though, only prefixes are available to construct solutions. Despite this limitation, our method performs comparably or even better in some cases, as shown in Figure 4. Since full-sequence prompts have more tokens, they show increased memorization in domains such as GitHub and books. To address this, we use a whitespace tokenizer to optimize prefixes, ensuring that performance remains competitive.

Victim as an Attacker LLM. We tested whether using the victim model as an attacker impacts performance and compared it with using distinct attacker models across different pre-training domains.

In prior experiments, the same model served as both the attacker and victim, but performance consistently lagged behind using Zephyr or GPT-4 as attackers. For example, with a sequence length of 200, Tulu LM as an attacker was 7.21% less effective than Zephyr, suggesting that using different attackers and sampling strategies significantly boosts performance.

Beyond GPT-4 for meta-prompt initialization.

In previous experiments, we employed GPT-4 for meta-prompt generation (see Section 3.2), but we now investigate the effect of using a less powerful open-source model on overall performance. Specifically, we utilize Mixtral-8x7B instruct (Jiang et al., 2024). In cases such as Alpaca with a sequence length of 200, Mixtral outperforms the prefix-suffix method, yielding 6.12% and 12.62% better reconstruction performance for base and instruct models, respectively, although it falls 4.00% short of GPT-4.

Training Data or Common Patterns. We test our method’s ability to generalize beyond pre-training data using the BookMIA dataset (Shi et al., 2023), which contains both training data members and non-members. Our method achieved a ROUGE-L score of 23.3 on training data members but only 16.7 on non-members, suggesting that our approach may lead the model to output memorized data rather than generalized information.

The impact of iteration count. Our method comprises two phases: sampling and refining. In the sampling phase, we use rejection sampling to gather data, and in the refining phase, we iterate three times on the most promising prompt, providing feedback at each step. Figure 5 illustrates performance improvements through these optimization stages. Although initial gains are modest from untargeted prompts, performance steadily improves

across iterations, peaking by the third round. Further iterations could enhance performance further but would come at higher computational costs.

8 Related Work

Data Extraction: Several studies have investigated data extraction techniques in LLMs. (Yu et al., 2023) proposed sampling adjustments for base models. (Nasr et al., 2023) focused on instruction-tuned models, demonstrating a divergence attack causing models like ChatGPT to repeat words indefinitely. (Zhang et al., 2023) developed a model interrogation attack to extract sensitive data by selecting lower-ranked output tokens. Additionally, (Geiping et al., 2024) introduced a system prompt repeater to extract sensitive system prompts, potentially compromising entire applications or secrets.

Jailbreaking: Emerging red-teaming methods exploit LLMs through jailbreaking techniques, aiming to coerce harmful behaviors (Shah et al., 2023; Li et al., 2023; Huang et al., 2023; Zeng et al., 2024; Mehrotra et al., 2023b; Hubinger et al., 2024). These approaches disrupt safety mechanisms, prioritizing harmful responses over data confidentiality.

9 Conclusion

In this work, we introduce a new method to analyze how instruction-tuned LLMs memorize pre-training data. Our empirical findings indicate that instruction-tuned models show higher memorization levels than their base models when using prompts that are different from the original pre-training data. However, this increased memorization in instruction-tuned models **does not imply** that these models regurgitate more data or are more vulnerable. Instead, it suggests that constructing instruction-based prompts reveals more pre-training data in instruction-tuned models.

Limitations

We would like to acknowledge that our method is mainly an auditing method which requires access to some part of the training data. We encourage future work to explore other automated strategies for building prompts for data extraction, targeting both base and instruction-tuned models, using prompts and contexts other than the original training data.

Ethics Statement

Enhancing the privacy-preserving capabilities of LLMs is crucial, given their increasing prominence and involvement in various aspects of life. Our new attack, designed to extract memorized data from instruction-tuned LLMs, which are widely used in real-world applications, deepens our understanding of these models' privacy limitations. By introducing this attack, we aim to advance the comprehension of memorization behaviors in different types of LLMs, encouraging future work to develop novel defense mechanisms to mitigate associated risks.

Acknowledgements

Funding support for the project activities of Niloo-far Mireshghallah and Yulia Tsvetkov has been provided by the Defense Advanced Research Projects Agency's (DARPA) SciFy program (Agreement No. HR00112520300), NSF DMS-2134012, IARPA HIATUS via 2022-22072200003, and ONR N00014-24-1-2207. The views expressed are those of the author and do not reflect the official policy or position of the Department of Defense or the U.S. Government. For Aly M. Kassem, this research is supported by the Vector Scholarship in Artificial Intelligence, provided through the Vector Institute and Natural Sciences and Engineering Research Council of Canada (NSERC) by NSERC Discovery Grant. This research was enabled in part by support provided by Compute Ontario and the Digital Research Alliance of Canada.

References

- Stella Biderman, USVSN Sai Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raff. 2023a. Emergent and predictable memorization in large language models. *arXiv preprint arXiv:2304.11158*.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023b. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. 2022. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*.

- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. 2021. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650.
- Nicolas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramer, Borja Balle, Daphne Ippolito, and Eric Wallace. 2023. Extracting training data from diffusion models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pages 5253–5270.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023a. [Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality](#).
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. 2023b. [Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality](#). See <https://vicuna.lmsys.org> (accessed 14 April 2023).
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Together Computer. 2023. [Redpajama: An open source recipe to reproduce llama training dataset](#).
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. 2023. [Enhancing chat language models by scaling high-quality instructional conversations](#). *Preprint*, arXiv:2305.14233.
- Michael Duan, Anshuman Suri, Niloofar Miresghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. 2024. Do membership inference attacks work on large language models? *arXiv preprint arXiv:2402.07841*.
- Ronen Eldan and Mark Russinovich. 2023. Who’s harry potter? approximate unlearning in llms. *arXiv preprint arXiv:2310.02238*.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.
- Jonas Geiping, Alex Stein, Manli Shu, Khalid Saifullah, Yuxin Wen, and Tom Goldstein. 2024. Coercing llms to do and reveal (almost) anything. *arXiv preprint arXiv:2402.14020*.
- Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muenighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, Will Smith, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah A. Smith, and Hannaneh Hajishirzi. 2024. Olmo: Accelerating the science of language models. *Preprint*.
- Michael M Grynbaum and Ryan Mac. 2023. The times sues openai and microsoft over ai use of copyrighted work. *The New York Times*, 27.
- Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2021. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing. *arXiv preprint arXiv:2111.09543*.
- Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. 2023. Catastrophic jailbreak of open-source llms via exploiting generation. *arXiv preprint arXiv:2310.06987*.
- Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera Latham, Daniel M Ziegler, Tim Maxwell, Newton Cheng, et al. 2024. Sleeper agents: Training deceptive llms that persist through safety training. *arXiv preprint arXiv:2401.05566*.
- Daphne Ippolito, Florian Tramèr, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher A Choquette-Choo, and Nicholas Carlini. 2023. Preventing generation of verbatim memorization in language models gives a false sense of privacy. In *Proceedings of the 16th International Natural Language Generation Conference*, pages 28–53. Association for Computational Linguistics.
- Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A Smith, Iz Beltagy, et al. 2023. Camels in a changing climate: Enhancing lm adaptation with tulu 2. *arXiv preprint arXiv:2311.10702*.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.
- Xuan Li, Zhanke Zhou, Jianing Zhu, Jiangchao Yao, Tongliang Liu, and Bo Han. 2023. Deepinception: Hypnotize large language model to be jailbreaker. *arXiv preprint arXiv:2311.03191*.

- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. 2023a. Tree of attacks: Jailbreaking black-box llms automatically.
- Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. 2023b. Tree of attacks: Jailbreaking black-box llms automatically. *arXiv preprint arXiv:2312.02119*.
- Fatemehsadat Mireshghallah, Archit Uniyal, Tianhao Wang, David Evans, and Taylor Berg-Kirkpatrick. 2022. [An empirical analysis of memorization in fine-tuned autoregressive language models](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1816–1826, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A Feder Cooper, Daphne Ippolito, Christopher A Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. 2023. Scalable extraction of training data from (production) language models. *arXiv preprint arXiv:2311.17035*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*.
- Jacob Pfau, Alex Infanger, Abhay Sheshadri, Ayush Panda, Julian Michael, and Curtis Huebner. 2023. Eliciting language model behaviors using reverse language models. In *Socially Responsible Language Modelling Research*.
- Aman Priyanshu, Supriti Vijay, Ayush Kumar, Rakshit Naidu, and Fatemehsadat Mireshghallah. 2023. Are chatbots ready for privacy-sensitive applications? an investigation into input regurgitation and prompt-induced sanitization. *arXiv preprint arXiv:2305.15008*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.
- Govind Ramesh, Yao Dou, and Wei Xu. 2024. Gpt-4 jailbreaks itself with near-perfect success using self-explanation. *arXiv preprint arXiv:2405.13077*.
- Avi Schwarzschild, Zhili Feng, Pratyush Maini, Zachary C Lipton, and J Zico Kolter. 2024. Rethinking llm memorization through the lens of adversarial compression. *arXiv preprint arXiv:2404.15146*.
- Rusheb Shah, Soroush Pour, Arush Tagade, Stephen Casper, Javier Rando, et al. 2023. Scalable and transferable black-box jailbreaks for language models via persona modulation. *arXiv preprint arXiv:2311.03348*.
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2023. Detecting pretraining data from large language models. *arXiv preprint arXiv:2310.16789*.
- Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Harsh Jha, Sachin Kumar, Li Lucy, Xinxu Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew E. Peters, Abhilasha Ravichander, Kyle Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Pete Walsh, Luke Zettlemoyer, Noah A. Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. 2024. Dolma: an Open Corpus of Three Trillion Tokens for Language Model Pretraining Research. *arXiv preprint*.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 2023a. Alpaca: A strong, replicable instruction-following model. *Stanford Center for Research on Foundation Models*. <https://crfm.stanford.edu/2023/03/13/alpaca.html>, 3(6):7.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023b. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, et al. 2023. Zephyr: Direct distillation of llm alignment. *arXiv preprint arXiv:2310.16944*.
- Tony Tong Wang, Adam Gleave, Tom Tseng, Kellin Pelline, Nora Belrose, Joseph Miller, Michael D Dennis,

- Yawen Duan, Viktor Pogrebnik, Sergey Levine, et al. 2023a. Adversarial policies beat superhuman go ais.
- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023b. [How far can camels go? exploring the state of instruction tuning on open resources](#). In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. *arXiv preprint arXiv:2304.01196*.
- Weichen Yu, Tianyu Pang, Qian Liu, Chao Du, Bingyi Kang, Yan Huang, Min Lin, and Shuicheng Yan. 2023. Bag of tricks for training data extraction from language models. *arXiv preprint arXiv:2302.04460*.
- Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyang Shi. 2024. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms. *arXiv preprint arXiv:2401.06373*.
- Yiming Zhang and Daphne Ippolito. 2023. Prompts should not be seen as secrets: Systematically measuring prompt extraction attack success. *arXiv preprint arXiv:2307.06865*.
- Zhuo Zhang, Guangyu Shen, Guanhong Tao, Siyuan Cheng, and Xiangyu Zhang. 2023. Make them spill the beans! coercive knowledge extraction from (production) llms. *arXiv preprint arXiv:2312.04782*.
- Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. 2024. [\(inthe\)wildchat: 570k chatGPT interaction logs in the wild](#). In *The Twelfth International Conference on Learning Representations*.
- 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 Hyperparameters Optimization

To ascertain the ideal hyperparameter balancing between memorization and overlap across diverse domains and sequence lengths, we initially streamlined our process by optimizing 20% of the dataset for quicker runtime. This entails iterating through multiple values to pinpoint the one that best aligns with our objectives. Subsequently, the selected values are applied to the entire dataset.

We select the following values for Llama-based models:

For a sequence length of 200, we allocate weights of 0.4 for memorization and 0.6 for overlap, a configuration tailored for C4, CC, and GitHub. Conversely, for ArXiv and Books, the emphasis shifts slightly, with 0.2 assigned to memorization and 0.8 to overlap.

At a sequence length of 300, nuances emerge across domains; for CC and C4, an even balance at 0.5 for memorization and overlap is determined. However, GitHub and ArXiv prefer a 0.4-0.6 split, favoring overlap slightly more. Conversely, Books lean towards a 0.3-0.7 ratio, emphasizing overlap more.

The weighting intensifies for a sequence length of 500, with C4, CC, and ArXiv converging at 0.5 for both memorization and overlap. GitHub adopts a 0.6-0.4 distribution, while Books adhere to a 0.4-0.6 allocation for memorization and overlap.

For the Falcon model, the designated values are as follows: For a sequence length of 200, we allocate a weight of 0.2 for memorization and 0.8 for overlap. With a sequence length of 300, the distribution shifts to 0.3 for memorization and 0.7 for overlap. Lastly, for a sequence length of 500, the weight is set at 0.8 for memorization and 0.2 for overlap.

B Detailed Results

B.1 Breakdown of Results from Section 5

In this section, we present a detailed breakdown of results for each instruction-tuned model, encompassing Alpaca, Tulu, and Vicuna, as depicted in Table 2. Figure 6 Shows a breakdown based on sequence length.



Figure 6: A detailed breakdown of the results presented in Table 1, over different sequence lengths and data domains for our proposed method. We can see that the instruction-tuned models demonstrate higher memorization scores (Rouge-L) compared to the base model. The full breakdown table, including the baseline methods, is provided in Appendix Table 2.

Alpaca-7B																	
Sequence	Method	Access	Github			ArXiv			CC			C4			Books		
			Mem	LCS _P	Dis	Mem	LCS _P	Dis	Mem	LCS _P	Dis	Mem	LCS _P	Dis	Mem	LCS _P	Dis
			↑	↓	↑	↑	↓	↑	↑	↓	↑	↑	↓	↑	↑	↓	↑
200	P-S-Base	B	.315	.125	-	.188	.107	-	.198	.103	-	.206	.111	-	.225	.090	-
	P-S-Inst	B	.294	.125	-	.200	.107	-	.168	.103	-	.152	.111	-	.153	.090	-
	Reverse-LM	B	.242	.191	.877	.141	.200	.890	.124	.203	.863	.117	.195	.862	.137	.151	.880
	GCG	W	.325	.107	.619	.189	.096	.473	.203	.087	.469	.214	.097	.404	.223	.077	.518
	Ours	B	.362	.102	.877	.205	.091	.890	.227	.101	.863	.213	.0939	.862	.247	.083	.880
300	P-S-Base	B	.295	.124	-	.186	.112	-	.193	.106	-	.208	.114	-	.213	.095	-
	P-S-Inst	B	.273	.124	-	.183	.112	-	.160	.106	-	.153	.114	-	.136	.095	-
	Reverse-LM	B	.232	.203	.881	.133	.145	.853	.117	.190	.822	.109	.182	.826	.123	.145	.877
	GCG	W	.311	.109	.535	.180	.100	.390	.197	.092	.378	.212	.102	.318	.200	.080	.432
	Ours	B	.330	.087	.881	.244	.110	.853	.222	.100	.822	.209	.094	.826	.228	.077	.877
500	P-S-Base	B	.263	.124	-	.175	.117	-	.179	.102	-	.196	.117	-	.184	.095	-
	P-S-Inst	B	.241	.124	-	.154	.117	-	.138	.102	-	.124	.117	-	.104	.095	-
	Reverse-LM	B	.214	.204	.833	.125	.192	.803	.099	.164	.805	.104	.167	.814	.105	.129	.838
	GCG	W	.265	.113	.435	.165	.107	.274	.182	.092	.274	.196	.113	.435	.173	.085	.317
	Ours	B	.275	.117	.833	.234	.122	.803	.193	.087	.805	.186	.083	.814	.189	.076	.838
Tulu-7B																	
200	P-S-Base	B	.315	.126	-	.188	.107	-	.198	.103	-	.206	.111	-	.225	.090	-
	P-S-Inst	B	.298	.125	-	.216	.107	-	.176	.103	-	.140	.111	-	.188	.090	-
	Reverse-LM	B	.254	.191	.877	.154	.200	.890	.130	.203	.863	.123	.195	.862	.153	.151	.880
	GCG	W	.325	.107	.619	.189	.096	.473	.203	.087	.469	.214	.097	.404	.223	.077	.518
	Ours	B	.372	.098	.877	.204	.093	.883	.225	.104	.858	.214	.095	.853	.236	.082	.882
300	P-S-Base	B	.315	.126	-	.188	.107	-	.198	.103	-	.206	.111	-	.225	.090	-
	P-S-Inst	B	.276	.124	-	.209	.112	-	.174	.106	-	.142	.114	-	.178	.095	-
	Reverse-LM	B	.246	.203	.881	.157	.196	.853	.125	.190	.822	.116	.182	.826	.134	.145	.877
	GCG	W	.311	.109	.535	.180	.100	.390	.197	.092	.378	.212	.102	.318	.200	.080	.432
	Ours	B	.341	.084	.878	.248	.108	.856	.222	.099	.824	.209	.090	.825	.231	.079	.872
500	P-S-Base	B	.263	.124	-	.175	.117	-	.179	.102	-	.196	.117	-	.184	.095	-
	P-S-Inst	B	.247	.124	-	.195	.117	-	.159	.102	-	.128	.117	-	.149	.095	-
	Reverse-LM	B	.233	.204	.833	.147	.192	.803	.107	.164	.805	.112	.167	.814	.118	.129	.838
	GCG	W	.265	.113	.435	.165	.107	.274	.182	.092	.274	.196	.113	.435	.173	.085	.317
	Ours	B	.363	.129	.814	.260	.112	.809	.216	0.079	.824	.207	.074	.829	.231	0.076	.841
Vicuna-7B																	
200	P-S-Base	B	.315	.126	-	.188	.107	-	.198	.103	-	.206	.111	-	.225	.090	-
	P-S-Inst	B	.311	.125	-	.225	.107	-	.215	.103	-	.205	.111	-	.212	.090	-
	Reverse-LM	B	.256	.191	.877	.199	.200	.890	.179	.203	.863	.180	.195	.862	.181	.151	.880
	GCG	W	.325	.107	.619	.189	.096	.473	.203	.087	.469	.214	.097	.404	.223	.077	.518
	Ours	B	.327	.094	.883	.199	.095	.888	.214	.100	.867	.200	.090	.866	.221	.083	.881
300	P-S-Base	B	.315	.126	-	.188	.107	-	.198	.103	-	.206	.111	-	.225	.090	-
	P-S-Inst	B	.267	.124	-	.194	.112	-	.208	.106	-	.182	.115	-	.189	.095	-
	Reverse-LM	B	.261	.203	.881	.204	.196	.853	.177	.190	.822	.173	.182	.826	.168	.145	.877
	GCG	W	.311	.109	.535	.180	.100	.390	.197	.092	.378	.212	.102	.318	.200	.080	.432
	Ours	B	.311	.078	.885	.241	.106	.854	.215	.097	.824	.201	.087	.833	.217	.076	.877
500	P-S-Base	B	.263	.124	-	.175	.117	-	.179	.102	-	.196	.117	-	.184	.095	-
	P-S-Inst	B	.241	.125	-	.219	.117	-	.193	.102	-	.188	.117	-	.192	.095	-
	Reverse-LM	B	.247	.204	.833	.198	.192	.803	.163	.164	.805	.166	.167	.814	.149	.129	.838
	GCG	W	.265	.113	.435	.165	.107	.274	.182	.092	.274	.196	.113	.435	.173	.085	.317
	Ours	B	.336	.116	.823	.255	.109	.817	.210	0.079	.823	.202	.075	.825	.233	0.078	.838

Table 2: Memorization scores (Mem), overlap between the prompts and suffix (LCS_P), and the distance between optimized and initial prompts (Dis) is evaluated across various pre-training data domains, evaluated across five scenarios: P-S-Base (sequence extraction on Llama), P-S-Inst (sequence extraction on the instruction-tuned model), Reverse-LM, GCG, and our method. Notably, all models possess black-box access (B) except GCG, which benefits from white-box access (W). The highest performance within each domain is highlighted in bold.

B.2 Improvement Percentages

To gauge the degree of enhancement relative to other baseline methods, we performed the following calculation: for each sequence length, domain, and model, we subtracted our method’s performance from that of each method and then divided the result by the performance of the other method. This allowed us to assess our method’s relative superiority or inferiority compared to the other method. The results shown in Table 3

Domain	Sequence Length	Alpaca			Tulu			Vicuna		
		P-S-INST	P-S-BASE	GCG	P-S-INST	P-S-BASE	GCG	P-S-INST	P-S-BASE	GCG
Github	200	.230	.149	.115	.249	.180	.145	.054	.039	.008
	300	.201	.119	.063	.232	.154	.096	.166	.055	.002
	500	.139	.042	.036	.467	.378	.370	.391	.273	.266
CC	200	.352	.144	.118	.279	.136	.111	-.003	.079	.055
	300	.387	.149	.127	.274	.146	.123	.030	.109	.087
	500	.399	.079	.062	.354	.206	.186	.089	.174	.156
C4	200	.401	.034	.005	.527	.035	-.004	-.022	-.029	-.066
	300	.367	.002	-.014	.469	.035	-.016	.107	-.034	-.051
	500	.497	-.005	-.053	.612	.057	.054	.075	.0297	.026
Books	200	.613	.095	.106	.250	.047	.057	.040	.018	-.009
	300	.681	.069	.142	.299	.081	.154	.144	.015	.084
	500	.809	.025	.089	.552	.252	.331	.210	.261	.340
ArXiv	200	.025	.090	.087	-.057	.080	.077	-.116	.057	.054
	300	.332	.313	.357	.187	.336	.380	.241	.296	.339
	500	.519	.334	.421	.331	.478	.574	.162	.449	.544

Table 3: Improvement percentages across diverse domains, sequence lengths, and models. P-S-INST denotes our method’s performance subtracted from P-S-INST performance and then divided on the latter, with similar comparisons for other methods.

B.3 Falcon Results

In this section, we present a detailed breakdown of results for the Falcon as depicted in Figure 7 with a breakdown based on sequence length.

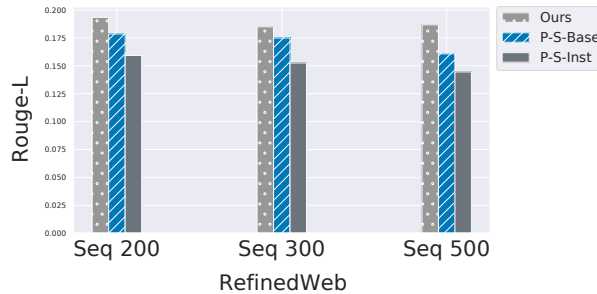


Figure 7: Comparison of our method to the P-S baseline on the Falcon model. We evaluate different sequence lengths of the pre-training data and observe that our method consistently outperforms the prefix-suffix base and instruction versions.

B.4 Common Patterns

To analyze the evolution from initial to optimized prompts, we examined common patterns by extracting the most frequent n-grams (n ranging from 1 to 5) in the optimized prompts. However, replacing these optimized n-grams with their counterparts in the initial prompts did not improve performance. This is because the transformation operates at the sentence level, where specific n-gram modifications—additions, deletions, or replacements—do not significantly impact the overall performance, given the complex interplay of various operations in the sentence-level transformation process.

B.5 Larger Sizes

In this section, we show the results for larger sizes, Alpaca-13B and Tulu-30B. We observed the same trend of our method in the larger sizes, as shown in Figure 8 and Figure 9. Note that we could only run

30B experiments on sequence length 200 and three subsets due to limited computational resources.

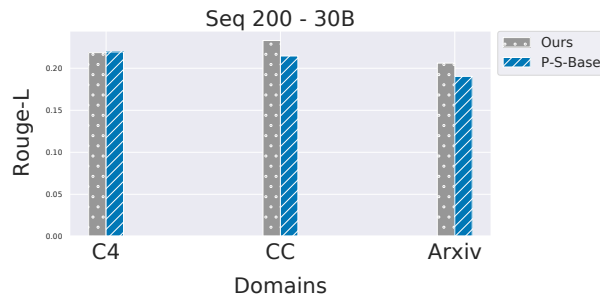


Figure 8: Comparison of our method to the P-S baseline on the Tulu-30B model. We evaluate different domains of the pre-training data and observe that our method consistently outperforms the prefix-suffix base and instruction versions.

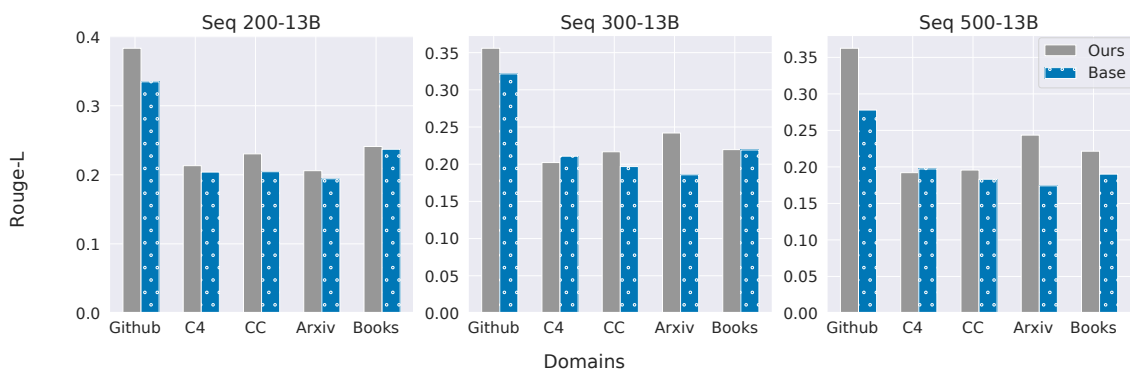


Figure 9: Comparison of our method to the P-S baseline on the Alpaca-13B model. We evaluate different domains of the pre-training data and observe that our method consistently outperforms the prefix-suffix base and instruction versions.

C Similarity Analysis on Different Instruction Tuned Models

This section delves into an error analysis of the instruction-tuned models utilizing the prefix-suffix and our optimization approach. We delve into the correlation, edit distance, and cosine similarity across the optimization prompt’s scores. Table 4 visually encapsulates the proximity of prompts from each model to one another. The initial part showcases the cosine similarity; notably, the similarity between the scores of the optimized prompts and the prefix-suffix exhibits lower similarity, while a substantially high similarity exists between the optimized prompts for each model, averaging around 90%.

Furthermore, upon computing the L_2 distance, a pattern emerges with a notable increase in distance between optimized prompts and prefix scores. Conversely, the distance shrinks significantly between the optimized prompts for various models. A similar trend unfolds in correlation analysis, wherein the correlation between the scores of the optimized prompts is notably high, contrasting with the lower correlation observed between the optimized and prefix-suffix.

These findings underscore the efficacy of the optimization process in generating very similar prompts for attacking various instruction-tuning models, which can indicate the universality of the optimized prompts.

<i>Cosine Similarity</i>					
Models (Ours)	Llama-7B (P-S-Base)	Tulu		Vicuna	
		P-S-Inst	Ours	P-S-Inst	Ours
Alpaca	.815	.835	.915	.838	.881
Vicuna	.822	.807	.903	-	-
Tulu	.837	-	-	-	-
<i>L₂-Distance</i>					
Alpaca	7.90	7.46	5.61	7.41	6.38
Vicuna	7.20	7.46	5.87	-	-
Tulu	7.50	-	-	-	-
<i>Correlation</i>					
Alpaca	.491	.512	.689	.477	.569
Vicuna	.410	.416	.636	-	-
Tulu	.509	-	-	-	-

Table 4: Comparison of Cosine Similarity, L2 Distance, and Correlation between Instruction-Tuned Models (Alpaca, Tulu, Vicuna) and Llama-7B using Prefix-Suffix and our proposed attack.

D Models & Evaluation Data Details

Attacker LLMs: Our attack strategy primarily relies on harnessing an open-source model known as Zephyr 7B β (Tunstall et al., 2023) as the attacker. This instruction-tuned variant of the Mistral-7B model has been fine-tuned on Ultra-Chat and Ultra-Feedback datasets (Ding et al., 2023) through DPO (Rafailov et al., 2024). Zephyr 7B β has demonstrated promising performance, particularly excelling in tasks related to writing and mathematics, despite its more compact size compared to larger models.

Victim LLMs We assess the memorization capabilities of instruction-tuned LLMs compared to their base model across various sizes by applying our attack on five open-source models of different sizes by employing the instruction-tuned versions of Llama (Touvron et al., 2023), OLMo (Groeneveld et al., 2024), and Falcon (Penedo et al., 2023). By comparing these instruction-tuned models to their base model, we gain insights into the impact of instruction-tuning on memorization.

Llama-based LLMs: Llama is known for its diverse instruction-tuned versions, each trained on various proprietary datasets. (1) Alpaca (7B, 13B; Taori et al. 2023a) is an early attempt at open-sourcing instruction-tuned models by fine-tuning on 52K instruction-following demonstrations generated from GPT-3.5. (2) Vicuna (7B Chiang et al. 2023b) is built through fine-tuning on 70K user-shared ChatGPT data, it showed competitive performance compared to OpenAI ChatGPT and surpassed Llama and Alpaca models. (3) Tulu (7B, 30B; Wang et al. 2023b) is fine-tuned on human+GPT data mixture of instruction-output pairs.

Falcon: The base model was trained on 1,000B tokens of RefinedWeb (RW) with curated corpora. We compare Falcon-Instruct 7B, an instruction-tuned version further trained on the Baize dataset (Xu et al., 2023).

OLMo: Open Language Models is a state-of-the-art 7 billion, open-source large language model released with full access to its inner workings and massive training data. OLMo trained on Dolma (Soldaini et al., 2024) with 2.5T tokens. We compare OLMo-Instruct 7B, an instruction-tuned version further trained on Tulu 2 SFT Mix and Ultrafeedback Cleaned (Iverson et al., 2023).

Data Domains To ensure comprehensive coverage of the pre-training data, we select 15,000 samples from five domains of the Llama data: Github (code), C4, CC (general knowledge), Arxiv (scientific papers), and Books. Each domain consists of 1,000 samples, totaling 5,000 for each of the three sequence lengths. For Falcon, we randomly select 3,000 samples from the RefinedWeb (RW), distributing 1,000 samples evenly across each sequence length. While for OLMo, we select 16,000 samples from six domains: The Stack (code), C4, CC (general knowledge), Reddit (social media), PeS2o (STEM papers), and Project Gutenberg (books). We followed the same splitting as in Llama, as each domain consists of 1,000 samples, totaling 6,000 for each of the three sequence lengths.

Sequence Lengths Selection To assess the resilience of our attack against different sequence lengths, we choose three: 200, 300, and 500. To better represent real-world usage, we choose the ratio of splitting each sample into prefix-suffix pairs based on analysis of the WildChat dataset (Zhao et al., 2024), which comprises 570K user-ChatGPT conversations spanning various languages and prompts. For each sequence length l , we provide the model with 33% of the sample as a prefix, while the remaining 67% serves as a suffix. For a length of 200 tokens, we allocate 66 for prefixes and 134 for suffixes. For 300 tokens, the divide is 100 for prefixes and 200 for suffixes. For 500 tokens, it is 167 for prefixes and 333 for suffixes.

GCG Inference Time It's worth noting that while GCG, which serves as the comparable baseline to our method, typically requires substantial resources and time to achieve convergence, our approach is significantly more efficient. Specifically, GCG takes approximately 12 minutes for a single sample to converge when running on two V100 GPUs. In stark contrast, our method completes the same task in just 1.30 minutes on the same hardware setup. This considerable computation time reduction highlights our approach's efficiency and effectiveness compared to the traditional GCG baseline.

E Examples of Instruction-Based Prompts

Prompt Type	Text	Mem ↑	LCS _P ↓
Initial Prompt	Generate a code snippet in Java that defines a class GetPrimaryKeysOperation which extends MetadataOperation. The class should be part of the package org.apache.hive.service.cli.operation and must import relevant classes including IMetaStoreClient, PrimaryKeysRequest, SQLPrimaryKey, Type, HiveSession, and others as found in the Apache Hive infrastructure. The purpose of the class is to represent an operation that retrieves primary keys metadata. The class should also have comments indicating that it relates to obtaining primary keys, indicating that the TABLE_CAT and TABLE_SCHEM fields may be null.	.096	.075
Optimized Prompt	How can one implement the GetPrimaryKeysOperation class in Apache Hive and what are the functions of IMetaStoreClient, PrimaryKeysRequest, SQLPrimaryKey, Type, and HiveSession during this process?	.490	.075

Prompt Type	Text	Mem \uparrow	LCS _P \downarrow
Initial Prompt	<p>Generate a code snippet for an XML layout used in an Android application that includes a vertical LinearLayout as the root element with match_parent attributes for both layout_width and layout_height. Within this LinearLayout, include two ImageView components, each with an id (main_imageview and main_imageview2 respectively), wrap_content for layout_width and layout_height. Also, add a Button with an id of testRx_button, match_parent for layout_width, and wrap_content for layout_height.</p>	.629	.562
Optimized Prompt	<p>Design a portrait orientation XML layout for an Android app using LinearLayout as the root element, with a fill_parent size attribute to occupy the entire screen horizontally and vertically? Inside this root LinearLayout, add two ImageView components side by side and place a Button component at the bottom, taking the entire horizontal width but a certain height. Make sure to give each element a unique ID during the process.</p>	.895	.125

Prompt Type	Text	Mem ↑	LCS _P ↓
Initial Prompt	<p>Generate a code snippet that includes a copyright and licensing header, common in many source files. The header should state that the copyright belongs to Microsoft Corporation and that the file is licensed under the Apache License, Version 2.0. Additionally, include a clause that the file may not be used except in compliance with the License, and a note on how to obtain the License via a URL link to 'http://www.apache.org/licenses/LICENSE-2.0'. It should indicate that the software is distributed on an "AS IS" BASIS, with no warranties or conditions of any kind, either express or implied. Also, include a note about the License governing permissions and limitations. After the header, write a using directive to include the 'System.Collections.Generic' namespace and declare a namespace called 'Microsoft.Azure.Commands.Compute'.</p>	.575	.523
Optimized Prompt	<p>In modern software development, its crucial to include a copyright and licensing header in your source files to clearly state the ownership and distribution rights of your code. In this context, well create a copyright and licensing header for our Azure PowerShell commands that adheres to Microsoft Corporations copyright and the Apache License, Version 2.0. Heres our code snippet:</p>	.938	.158

Prompt Type	Text	Mem ↑	LCS _P ↓
Initial Prompt	Can you create an instructional guide on the topic "How to Remove my FB Account," addressing concerns such as feeling overwhelmed by unnecessary posts, desiring privacy, and annoyance from continuous notifications, while also providing a solution for those undecided on whether to leave permanently by explaining the difference between deactivating (which maintains friendships for potential future reactivation) and permanently deleting a Facebook account, and ends with an instruction to 'click' for the final step?	.148	.095
Optimized Prompt	permanently or temporarily remove my presence from Facebook to avoid irrelevant content, notifications, and protect my privacy? What are the variances between these two options?	.269	.043

Prompt Type	Text	Mem \uparrow	LCS _P \downarrow
Initial Prompt	What is a quilting technique that involves creating patterns such as letters and numbers by maneuvering the fabric freely under the needle, using colored threads and flowing letter stencils to enhance the design, which also involves tracing the design onto lightweight paper, pinning to the fabric, and sewing along the lines with the feed dogs disengaged as per the sewing machine's owner's manual?	.049	.028
Optimized Prompt	use free motion quilting to create text designs using colored threads and letter stencils? What is the process involved, including transferring the design onto paper, securing it to the fabric, and sewing with the feed dogs disengaged, as detailed in the sewing machine manual?	.288	.125
Initial Prompt	What is the content of the Introduction section, labeled 'sec1,' that outlines the origins of the directed power graph notation $\vec{P}(S)$ of a semigroup S , as established by Kelarev and Quinn, and includes the definition provided by these authors in which each arc represents an exponentiation relationship between semigroup elements, as well as the subsequent definition of an (undirected) power graph $P(S)$ by Chakrabarty et al., along with its criterion for vertex adjacency?	.236	.253
Optimized Prompt	In the works of Kelarev and Quinn, as well as in the research by Chakrabarty et al., what is the significance behind the notation $\vec{P}(S)$ for directed power graphs, and how does it differ from the undirected version $P(S)$ that they all define?	.400	.106

Prompt Type	Text	Mem \uparrow	LCS _P \downarrow
Initial Prompt	Can you create an introductory paragraph for a mathematical text that defines the exponential growth rate of a finitely generated group with respect to a finite generating set, detailing the set of elements within a given word length as well as the formula used to determine whether the group has exponential growth based on the limit of the cardinality of that set to the power of the reciprocal of the word length?	.195	.169
Optimized Prompt	How can we understand the concept of exponential growth rate in the study of finite groups, specifically in terms of the size of sets of elements with a fixed word length and a formula based on the limit of these sizes raised to the power of the word lengths reciprocal? This section will define this growth rate and elucidate its importance in the context of group theory.	.366	.112

Prompt Type	Text	Mem ↑	LCS_P ↓
Initial Prompt	What are the key differences between Certificates of Deposits (CDs) and government bonds as investment options according to MyBankTracker, and how does the explanation by Simon Zhen help an individual with limited resources determine which investment is more suitable for their savings strategy?	.185	.202
Optimized Prompt	How does MyBankTracker differentiate between Certificates of Deposit (CDs) and government bonds, and how can someone with limited resources determine which investment option is more suitable for their savings strategy based on Simon Zhens explanation?	.292	.080

Prompt Type	Text	Mem ↑	LCS _P ↓
Initial Prompt	<p>Can you provide an account of the narrative presented on "This American Life" about the incident from the summer of 1951 in small-town Wisconsin, where two baby girls were accidentally switched at birth and taken home by the wrong families, focusing on how host Ira Glass introduced the characters Kay McDonald and Mary Miller, the impact of Mary Miller revealing the secret after 43 years through letters to Sue and Marti, the daughters involved, and the exploration of the emotional aftermath by reporter Jake Halpern, including the perspectives of the mothers and their struggle with the truth, as part of an episode which also featured other segments such as a historical article about a slave auction, a review of William Kane's case, and a segment titled "Strength In Numbers"?</p>	.126	.219
Optimized Prompt	<p>Could you retell the tale shared on This American Lives podcast from the summer of 1951 in a small Wisconsin town, detailing the unintentional swapping of newborns between families bearing the names Kay McDonald and Mary Miller? Please include the introduction of critical characters, the ramifications brought about by Mary Millers disclosure following forty-three years, as well as the sentimental reaction explored by reporter Jake Halpern, while also mentioning any other sections included in the episode.</p>	.241	.103