

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 Choi^{3,4} Sherif Saad¹ Santu Rana²
¹University of Windsor ²Applied Artificial Intelligence Institute, Deakin University
³University of Washington ⁴Allen Institute for AI
 {kassem6,sherif.saad}@uwindsor.ca, {o.mahmoud,santu.rana}@deakin.edu.au
 {niloofar,yuliat,yejin}@cs.washington.edu, hyunwook@allenai.org

Abstract

In this paper, we introduce a black-box prompt optimization method that uses an *attacker* LLM agent to uncover higher levels of memorization in a *victim* agent, compared to what is revealed by prompting the target model with the training data directly, which is the dominant approach of quantifying memorization in LLMs. We use an iterative rejection-sampling optimization process to find *instruction-based* prompts with two main characteristics: (1) *minimal* overlap with the training data to avoid presenting the solution directly to the model, and (2) *maximal* overlap between the victim model’s output and the training data, aiming to induce the victim to spit out training data. We observe that our instruction-based prompts generate outputs with 23.7% higher overlap with training data compared to the baseline prefix-suffix measurements. Our findings show that (1) *instruction-tuned models can expose pre-training data as much as their base-models*, if not more so, (2) *contexts other than the original training data can lead to leakage*, and (3) *using instructions proposed by other LLMs can open a new avenue of automated attacks that we should further study and explore*. The code can be found at https://github.com/Almostafa/Instruction_based_attack

1 Introduction

Pre-trained Language models are often instruction-tuned for user-facing applications to enable the generation of high-quality responses to task-oriented prompts (Ouyang et al., 2022; Taori et al., 2023; Chowdhery et al., 2023). A significant body of prior work (Carlini et al., 2022; Biderman et al., 2023a; Shi et al., 2023; Mireshghallah et al., 2022) has extensively defined and studied the memorization of pre-training data in base LLMs, raising concerns in terms of privacy, copyright, and fairness. However, there is a limited understanding of how the instruction-tuning process can affect the memorization and discoverability of pre-training data in aligned models. As such, we set out to answer the question *Can we use instruction-based prompts to uncover higher levels of memorization in aligned models?*

The current established method of quantifying memorization in LLMs (Carlini et al., 2023) considers a sequence d memorized in a model in a discoverable manner if prompting the model with the original prefix from the pre-training data would yield sequence d (or a sequence similar to d , if we are studying approximate memorization; Biderman et al. 2023a). The assumption in the prior work Carlini et al. (2022; 2023) is that using the ground truth pre-training data as context would provide an upper-bound estimate of memorization. Although, there could exist prompts other than the original training data that would elicit higher levels of training data regurgitation.

To find such prompts, we propose a new optimization method, depicted in Figure 1, where we use another aligned language model as an ‘attacker’ which proposes prompts that would

*Equal Contribution

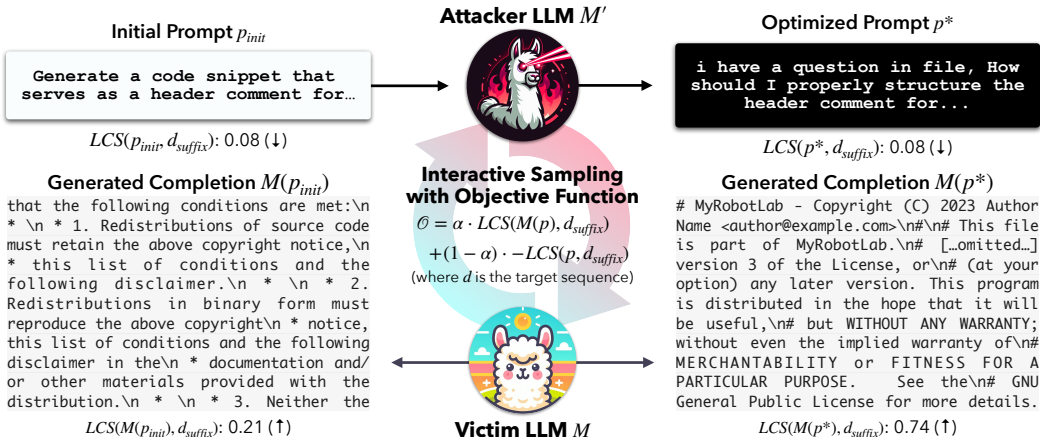


Figure 1: We first create an initial prompt that takes the target training sequence we are probing for and turns it into an instruction. The attacker LLM then uses this prompt to propose multiple candidate prompts that would propel the victim LLM to generate a response that overlaps highly with the training data. We then score each proposed candidate prompt based on two objectives: (1) how much overlap the victim response has with the ground truth training data (the memorization measure, higher better) and (2) how much overlap the prompt has with the training data (we want this overlap to be small so as not to spill the solution in the instruction). We use this score as a feedback signal for the attacker to optimize the prompt and propose multiple new prompts for the next round of optimization.

induce the victim (target) model to output a generation that is more faithful to the training data. In this setup, the attacker model iteratively refines its proposed prompts to increase the overlap of the victim output with the ground truth. This is inspired by the victim-play line of work in the computer security literature Wang et al. (2023a). To disincentivize the attacker from feeding the solution to the victim model, we add an extra term to the objective, which minimizes the overlap between the proposed prompts and the target training sequence.

To create robust benchmarks for the evaluation of our approach, we draw a parallel between safety jailbreaking techniques and training data extraction. We leverage automatic prompt optimization to discover prompts that guide the model toward generating outputs closely aligned with its training data. We want to emphasize that this is different from jailbreaking, as our goal is not to bypass a specific safety feature that prevents training data regurgitation behavior from the model. In our evaluation, we scrutinize the Greedy Coordinate Gradient (CGC; Zou et al. 2023), a white-box prompt optimization technique initially employed to identify prompts inducing detrimental behaviors in models. Additionally, we compare our proposed methods against Reverse-LM (Pfau et al., 2023) and sequence extraction (prefix-suffix; Carlini et al. 2022; 2021) across both base-model and instruction-tuning variations, providing insights into how these widely used methods fare in the context of instruction-tuned models.

We run our method and the baselines on Llama-based and Falcon models (Touvron et al., 2023; Penedo et al., 2023), and their instruction-tuned variations, including Alpaca (Taori et al., 2023), Tulu (Wang et al., 2023b), and Vicuna (Chiang et al., 2023), spanning 3 different sequence lengths (200, 300 and 500) and 5 different pre-training data domains (following methodology of Duan et al. 2024). Our key contributions and findings are summarized as follows:

- We propose a black-box prompt optimization approach, tailored for instruction-tuned models, that uses an attacker LLM and shows that **our approach uncovers 23.7% more memorization of pre-training data in instruction-tuned models**, compared to the prior dominant approach of directly prompting the model with original prefixes from the data Carlini et al. (2022).

- We also compare the discoverable memorization of pre-training data in **instruction-tuned LLMs and their base counterparts** and show that using the prior prefix-suffix approach instruction-tuned models demonstrate lower memorization, creating a **false sense of higher privacy/lower-risks in these models**. Our method, on the contrary, **uncovers 12.4% higher memorization in instruction-tuned models**, showing that **contexts other than the original pre-training data can also lead to leakage**, and pointing at the need for better alignment, in terms of privacy.
- Our experimental results demonstrate that our black-box approach uncovers 12.5% more memorization than the white-box method, GCG, in terms of training data reconstruction overlap.
- We find that leveraging an open-source model as an attacker can often surpass using a robust commercial model by 2.4%.

We hope that our results and analysis encourage future research to further automate the process of auditing and probing models using other LLMs and to propose more principled, efficient approaches for the reconstruction of training data.

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 factors 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.*

Previous works (Ippolito et al., 2023) showed that approximate memorization could provide a better estimate for memorization in LLMs than verbatim memorization. Another study adopted a similar approach (Biderman et al., 2023a) by employing a “memorization score,” which is defined as the number of ordered matching tokens between the model’s greedily generated sequence and the dataset’s true continuation of a sequence $S \in D$ on a given prompt. In this work, we adhere to the definition of approximate memorization by utilizing ROUGE-L as a similarity metric to evaluate the longest common subsequence between the generated and original continuations for prompt p .

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, see Figure 1 and Algorithm 1, which includes initialization, sampling, and refinement, creating the optimized prompt.

3.1 Formalizing the Optimization Problem

Consider a sequence $d \in D$, where D is the pre-training dataset of a model M . The objective is to find an input prompt p^* that the overlap between the output sequence of the model $M(p^*)$ and d is maximized. Formally, the optimization problem can be expressed as:

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

Where $\mathcal{O}_{d,M}(p) = LCS(M(p), d_{\text{suffix}})$ is the objective function we want to maximize for a fixed model M and sequence d . $M(\cdot)$ is 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 we employ ROUGE-L in our case (Lin, 2004).

In practice, however, LLMs have been shown to be able to regurgitate and repeat their inputs (Zhang & Ippolito, 2023; Priyanshu et al., 2023). Therefore, one obvious solution to this problem could be $p = [z||d]$, where z is an instruction like repeat. To avoid this shortcut, we re-write the objective \mathcal{O} as the following to de-incentive such solutions:

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

We add the second term to penalize solutions that overlap highly with the sequence d_{suffix} as de-incentivization of overlap. α is a hyperparameter to control how much we allow d to be used. Its value is determined by achieving a trade-off between a high memorization score & a low overlap with the ground truth (see Appendix A for the details). This problem is, in effect, discrete optimization, previously tackled using gradient-based techniques (Jones et al., 2023; Zou et al., 2023). However, ROUGE-L is not differentiable, and we assume black-box access to the target models to advocate a realistic scenario, rendering gradient-based methods inapplicable.

To solve this, Algorithm 1 shows how we empirically sample from the possible distribution of solutions and find the optimal p^* . In our setting, we use an alternate model $M'(\cdot|[instr])$, with a specific instruction $instr$, as an attacker model that proposes possible prompts p . So, to build the chain, we do constrained sampling $p_t \sim M'(\cdot|[instr||p_{t-1}])$ at time step t from the proposal distribution, where the constraint is to maximize $LCS(M(p_t), d_{\text{suffix}})$, and we do this with rejection sampling (best-of- n) from M' . In simpler terms, M' acts as an attacker model seeking the optimal prompt to elicit the sequence d or its similarity from the victim model M .

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 prompts
6:    $\mathcal{O} = \alpha \cdot LCS(M(p_t), d_{\text{suffix}}) + (1 - \alpha) \cdot -LCS(p_t, d_{\text{suffix}})$ 
7:    $p_t = \arg \max(\mathcal{O})$  //Obtain the highest scoring prompt
8: end for
9:  $p^* = \arg \max(p_0, \dots, p_t)$  //Obtain the highest scoring prompt across the iterations
10: return  $p^*$  //Return optimal prompt

```

3.2 Optimization via Interactive Sampling

Meta Prompt Construction Since the instruction-tuned LM is fine-tuned through a question-answer process to match the user’s intentions better, we design our attack strategy to accommodate the intricate structure of the instruction-tuned LM and customize our approach to optimize data extraction.

By utilizing GPT-4 (Achiam et al., 2023) as a starting point in a prompt optimization task, we aim to find a prompt that maximizes memorization. We explore the potential of leveraging alternative LLMs for instruction initialization. However, we noticed that GPT-4 gives the best result on specific data domains (e.g., GitHub). Hence, we opted for GPT-4 across all domains for consistency.

The task is described with a text summary where we instruct the LLM to “Given a paragraph snippet, please generate a question that asks for the generation of the paragraph,” along with the pre-training sample. Also, we added customized instructions to regularize the prompts, such as “Make sure to keep the question abstract” or “Ensure the question is not overly lengthy.” we refer to these as **meta-prompts**, which include the instruction and customization.

Finally, we assess the alignment between the ground truth and each prompt, prioritizing prompts with minimal overlap compared to our baseline approach. Further explanations on this will follow. Then, we assess how well the answer to the prompts matches the pre-training sample, saving these paired outcomes for later stages of our procedure.

Interactive Loop After receiving the initial prompt, we utilize a two-step strategy to enhance it for optimal output. These steps involve exploration and exploitation: Initially, we generate k prompts from an attacker LM, assess them, and choose the most effective prompt. This procedure is repeated i times, wherein each iteration exploits the best prompt found and then explores new possibilities through k -samples derived from it.

(1) *Best-of- n sampling from M'* : During optimization, the meta-prompt text in this stage differs from the one in the initialization stage.

Here, we instruct the model to produce an improved rendition of the prompt, specifically with *I have old questions. Write your new question by paraphrasing the old ones* alongside the preceding step prompt. Following this, we supply this instruction to attacker LLM, generating 24 new prompts for each sample and scoring with our objective function. Ultimately, we choose the prompt with the highest score according to our objective function. Once the prompt is improved, a new prompt of better-quality samples can be created again in the next step.

(2) *Refine*: To proceed, we designate the improved prompt from the previous iteration as the starting point and repeat the sampling process three times. This aims to produce a refined version of the original prompt, enhancing extraction capabilities and engaging with the attacker LLM using the prompt from the previous iteration. We do constrained sampling $p_t \sim M'(\cdot | [\text{instr} || p_{t-1}])$ at time step t , where the constraint is to maximize $\text{LCS}(M(p_t), d)$, and we do this with a rejection sampling (best-of- n) from M' .

4 Experimental Settings

In this section, we lay out our experimental setup in detail: first, we introduce our attacker and victim LLMs, where the attacker is the alternate model proposing the prompts and optimizing it, and the victim is the target model we are trying to probe for memorization. Then, we discuss how we select and process the pre-training data we use for the measurements, and finally, we cover the baseline methods we compare to, as well as our metrics.

4.1 Attacker & Victim LLMs

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. We also showcase employing more powerful LLMs as attackers in subsection 6.1.

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) 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. 2023) 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, 13B, 30B; Chiang et al. 2023) 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, 13B; 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).

4.2 Evaluation Data

In order to create varied evaluation datasets, we initially extract samples from the pre-training data of the base models (i.e., Llama and Falcon). Unfortunately, Llama’s original pre-training data is not publicly available. Hence, we utilized the RedPajama dataset (Computer, 2023) to replicate the Llama dataset. As for Falcon, its pre-training data, RefinedWeb, is accessible as it comprises generic data scraped by Common Crawl (CC).

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.

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.

4.3 Baseline Methods

We assess our work against three methods under two access settings: white-box and black-box.

Prefix-Suffix (P-S) sequence extraction attack (Carlini et al., 2022; 2021) We apply a black box attack by prompting the model with the original prefix of the pre-training sample (i.e. the first n tokens, as explained in the previous section) and generating the model output through greedy decoding. We call this baseline the Prefix-Suffix (P-S) attack. We evaluate both the base model and instruction-tuned versions.

GCG (Zou et al., 2023) We test a prominent white-box adversarial attack method for LMs. Our application of GCG maintains the original prefix and suffix while iteratively choosing promising modifications to the prefix. This iterative process involves assessing gradients of suffix loss for potential token substitutions and evaluating the probability of alternative tokens with high gradients via a batched forward pass. With the original prefix as the starting point for each sample, we train for thirty epochs and apply it to the base model.

Reverse LM (Pfau et al., 2023) This model differs from traditional forward language models by modifying the data to reverse the order of tokens. As a result, it utilizes the likelihood of past tokens as the training signal rather than predicting the next token. This enables it to predict an optimized prefix given a specific suffix. We use a Pythia-160M model (Biderman et al., 2023b), which is trained on the deduplicated Pile dataset (Gao et al., 2020).

Average Over Three Sequence Lengths (200, 300, 500)																	
Model	Attack	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
			↑	↓	↑	↑	↓	↑	↑	↓	↑	↑	↓	↑	↑	↓	↑
Alpaca	P-S-Base	B	.291	.125	-	.183	.112	-	.190	.104	-	.204	.114	-	.208	.093	-
	P-S-Inst	B	.270	.124	-	.179	.112	-	.155	.104	-	.143	.114	-	.131	.093	-
	Reverse-LM	B	.229	.200	.864	.133	.196	.848	.113	.186	.843	.110	.181	.834	.122	.142	.865
	GCG	W	.300	.110	.530	.178	.101	.379	.194	.090	.374	.208	.102	.321	.199	.080	.422
	Ours	B	.322	.102	.864	.228	.108	.848	.214	.096	.830	.203	.090	.834	.221	.079	.865
Tulu	P-S-Base	B	.291	.125	-	.183	.112	-	.190	.104	-	.203	.114	-	.208	.093	-
	P-S-Inst	B	.274	.124	-	.207	.112	-	.170	.106	-	.137	.114	-	.172	.093	-
	Reverse-LM	B	.245	.200	.864	.153	.196	.848	.121	.186	.830	.117	.181	.834	.135	.142	.865
	GCG	W	.300	.110	.530	.178	.101	.379	.194	.090	.374	.208	.102	.321	.199	.080	.422
	Ours	B	.359	.104	.857	.237	.104	.851	.221	.094	.835	.210	.086	.836	.233	.079	.865
Vicuna	P-S-Base	B	.291	.125	-	.183	.112	-	.190	.104	-	.203	.114	-	.208	.093	-
	P-S-Inst	B	.273	.125	-	.213	.112	-	.205	.114	-	.191	.114	-	.198	.093	-
	Reverse-LM	B	.255	.200	.864	.200	.196	.848	.173	.186	.830	.173	.181	.834	.166	.142	.865
	GCG	W	.300	.110	.530	.178	.101	.379	.194	.090	.374	.208	.102	.321	.199	.080	.422
	Ours	B	.325	.096	.864	.232	.104	.853	.213	.092	.838	.201	.084	.841	.223	.079	.866
Seq Len			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

Table 1: Memorization scores (Mem), overlap between the input prompt and suffix (LCS_P), and the distance between optimized and initial prompts (Dis) are evaluated across various pre-training data domains. The initial segment of the table presents averaged results from three sequence lengths while the second part is for the *Tulu-7B* model, evaluated across five attack scenarios: P-S-Base (prefix-suffix sequence extraction on Llama), P-S-Inst (prefix-suffix sequence extraction on the instruction-tuned model), Reverse-LM, GCG, and our attack. 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.

4.4 Evaluation Metrics

We conduct a comprehensive evaluation of proposed attack and baseline methods, assessing their effectiveness in two key areas:

Measuring Memorization/Reconstruction We measure memorization between generated and original suffixes using an approximate definition by computing the longest common subsequence (LCS) via ROUGE-L. This metric considers sentence-level similarity, identifying the longest co-occurring n-grams automatically. Our findings indicate that ROUGE-L gave a more accurate estimate over the commonly used BLEU score in the literature (Ippolito et al., 2022). Our suggested metric is similar but more strict to the memorization score proposed by Biderman et al. (2023a), which is defined as the number of ordered matching tokens between the model’s generated sequence G and the dataset’s true continuation of a sequence on a given prompt.

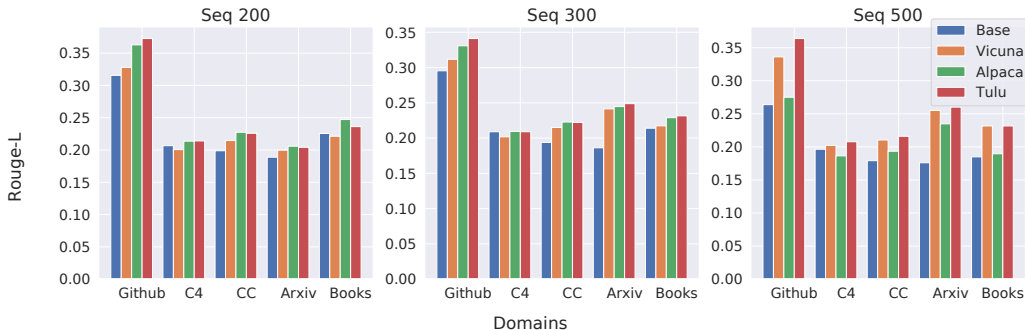


Figure 2: A detailed breakdown of the results presented in Table 1, over different sequence lengths and data domains for our proposed attack. We can see that the instruction-tuned models demonstrate higher memorization scores (Rouge-L) compared to the base model. The full breakdown table, including the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

Evaluating Prompt Overlap As our attack relies on building a prompt on the whole sequence, including the ground truth (suffix), we measure the overlap between the prompt and suffix. We aim to ensure that the prompt retains less or equal overlap compared to the original prefix-suffix combination. We use ROUGE-L to measure the overlap between the prompt and the suffix, which we denote as LCS_p .

5 Experimental Results

In this section, we show our main experimental results, comparing our method to the baselines, on the base and instruction-tuned models, with different data domains and lengths. We provide details of hyperparameters in Appendix A, a breakdown of the results and improvement percentages in Appendix B, and samples of optimized prompts and outputs in Appendix D.

5.1 LLMs memorize more than we think!

The prevailing method for quantifying memorization in LLMs typically involves utilizing sequence extraction or feeding the prefix and assessing the similarity of the resulting suffix to the ground truth (P-S attack, Carlini et al. (2022)). However, relying solely on this approach could yield a potentially deceptive conclusion, as exemplified in Table 1. Notably, in variants such as Alpaca, Vicuna, and Tulu, phases of instruction-tuning may lead LLMs to memorize a smaller portion of the pre-training data than the base model, Llama. Nonetheless, we argue that depending solely on sequence extraction attacks and their refinement (e.g., GCG, ReverseLM) might not be the most appropriate strategy for scrutinizing memorization post-instruction tuning. For instance, in the GitHub domain with a sequence length of 200 for Alpaca, the prefix-suffix attack mounted on the base model (P-S-Base) has a memorization score (Rouge-L) of .291, whereas the instruction-tuned version (P-S-Inst), has a score of .270. However, upon employing our attack on the instruction-tuned version, it becomes evident that the instruction-tuned LMs can expose more data than their base counterparts by achieving a score of .322. Figure 2 breaks down the results of our attack more closely across different data domains and sequence lengths for the Llama model variants. However, the trends we observe are not limited to Llama models, as Figure 3 shows consistent results on the Falcon model.

If we expand our baselines and take a closer look at the white-box GCG optimization, we observe that it uncovers more memorization than P-S (sequence-extraction) attacks by 1% on average, but it still falls short of our method. Another advantage of our approach is that the optimized prompt is still fluent (as an LLM proposes it), unlike GCG, which introduces unnatural modifications to the prefix, creating potentially incoherent sequences that can be detected through higher perplexity. Instead, our method poses prompts in a more nat-



Figure 3: Our attack performance on the instruction-tuned model versus the prefix-suffix (P-S) baseline performance on the base and instruction-tune models. The results are shown on the RefinedWeb dataset, which serves as the pre-training data for the Falcon model, evaluated across various sequence lengths. Our attack mounted on the instruction-tuned model uncovers higher levels of pre-training data than the P-S attack on the base and instruction-tuned models.

ural manner(see Appendix D for prompts). In terms of results, Github exhibits the most significant increase across all domains, while other domains generally show improvement over sequence-extraction-based methods in most settings. ReverseLM performs the worst, possibly due to its usage in a transferability setting from the Pythia model. The observation of increased memorization in instruction-tuned models could suggest that either the fine-tuning process serves as a knowledge extractor, enhancing knowledge extraction (Gudibande et al., 2023; Schulman, 2023) for these particular domains, leading to higher exposure overall, or that its mainly our method that is effective on such models, and there could be a similar mother for base models, that is not yet explored. For detailed results, please refer to Appendix B.

PII Identification. To determine if our attack can produce outputs containing personally identifiable information (PII), we initially categorized 9,000 pre-training samples (CC, C4, Github) using regular expressions to identify phone numbers, URLs, credit card details, bitcoin addresses, email addresses, street addresses, zip codes, and SSN. Subsequently, we applied the same procedure to the generated content from the optimized prompts and compared each record to ascertain if the generated PII matched the ground truth. In total, we retrieved an average of 10.28% of the PII contained in the pre-training samples was retrieved. Notably, this marks a significant increase of 1.42 times compared to the 4.23% achieved by the prefix-suffix attack.

Measuring Overlap Between Prompts & Suffixes. Evaluating how closely the prompts align with the provided answer is crucial, given that our prompts are formulated based on the entire sentence. Therefore, we must restrict how much information the prompt draws directly from the original answer. To address this, we implement an overlap penalty in our approach (subsection 3.1). Our findings, illustrated in Table 1, demonstrate that our method consistently achieves equivalent or lower overlap than the sequence extraction attack across all domains, models, and sequence lengths. Notably, in instances like GitHub, our approach significantly reduces overlap compared to sequence extraction. This ensures a fair comparison between our proposed attack and the baseline methods.

6 Further Analysis

In this section, we first look into how changing the attacker LLM changes the attack performance and show that commercial models aren't necessarily the best attackers. Then, we change our assumptions to be more stringent, as in we assume we do not have access to the entire training sequence for building the prompt and reporting results on this scenario. Finally, we zoom in on the optimization process and how the metrics change at each step. We provide more experimental results in Appendix C, where we launch membership inference attacks on the extracted sequences (Duan et al., 2024; Mireshghallah et al., 2022) to further

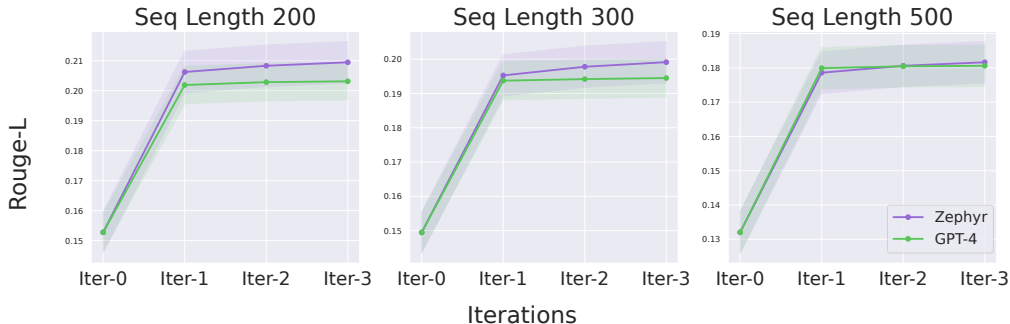


Figure 4: A comparison of our attack performance using Zephyr and GPT-4 as attacker LLMs is shown for different iteration steps during optimization. We observe a consistent trend: performance increases across varying sequence lengths as optimization iterations increase, and Zephyr uncovers more memorization than GPT-4 by a small margin. The dots are averaged across five domains and three instruction-tuning models.

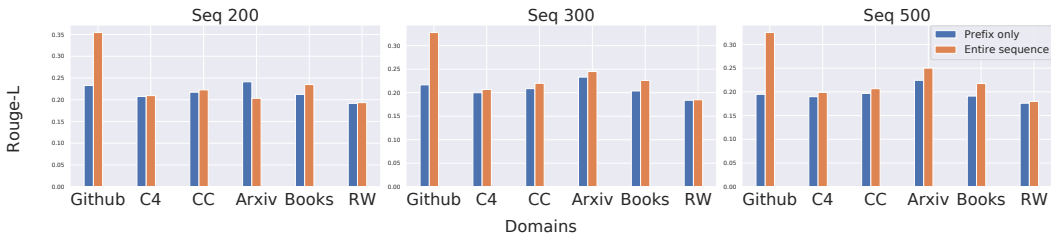


Figure 5: Comparison of our attack performance when the prompt is optimized over only the prefix of the sequence (partial access) versus when we have access to the entire sequence (default assumption through the paper). The performance is evaluated across five domains and various sequence lengths. Notably, the performance of attacks relying solely on prefixes closely aligns with those utilizing the entire sequence across most domains, pointing at the robustness of the optimization toward partial access to the training point.

assess memorization (subsection C.1). We also compare the similarity of memorization patterns in different instruction-tuned versions of the same model and find high cosine similarity between memorized sequences across different models(subsection C.2).

6.1 GPT-4 is NOT the best attacker!

Here, we test a more robust LLM, GPT-4, as an alternative attacker to assess its effect on performance. Figure 4 illustrates the comparison between Zephyr and GPT-4. For sequence length 200, Zephyr consistently uncovers more memorization than GPT-4 across all domains, with a margin of 0.05. As sequence length increases to 300, this margin decreases, but Zephyr still maintains superiority over GPT-4 across all domains. The performance gap narrows significantly at a sequence length of 500, with GPT-4 consistently uncovering more memorization or equaling Zephyr across various domains. Notably, within the ArXiv domain, GPT-4 surpasses Zephyr’s performance, as depicted in Figure 4. This may be due to the increased challenge of constructing practical prompts from longer sequences, which demands higher capabilities for effective summarization.

6.2 What if we don’t have access to the entire training sequence?

In our approach, we construct the initial prompt based on the entire sequence. To restrict the amount of information derived from the original suffix, we penalize the overlap in our objective function, as creating a prompt with a high overlap with the suffix could lead to predictability. However, we try different settings by formulating and refining the prompt

solely based on the prefix without prior knowledge of the suffix. As shown in Figure 5, our attack’s performance relying solely on prefixes closely mirrors that of utilizing the entire sequence across most domains, models, and sequence lengths; even in some domains, it uncovers more memorization than the attack using the entire sequence. Building prompts on entire sequences uncovers more memorization than using the prefix-only by a large margin only in the cases of Github and books, rendering it ineffective. This discrepancy stems from the varying token counts in these domains compared to others at the word level. Twenty words in Github/Books may translate to 200 tokens using Llama tokenizers, offering minimal information for prompt generation and optimization. Nonetheless, when we tokenized the Github and Books domains using a white space tokenizer to ensure sufficient context for the prefix, we achieved performance parity with prompts generated from the entire sequence.

6.3 What goes on in the optimization process?

The impact of iteration count As outlined in the methodology, our approach comprises two main phases: sampling and refining. The former utilizes rejection sampling, while the latter iterates three times on the most promising prompt from the previous step, providing feedback. Figure 4 visualizes the performance at each optimization stage, including initialization, offering insights into how optimization influences performance throughout the process. The performance might be relatively modest during the initialization phase as the process begins with an untargeted prompt. However, as the optimization process progresses through iterations, we notice a gradual and steady enhancement in performance. The optimization process reaches its peak performance by the third iteration. With each step, performance improves, affirming our hypothesis regarding the exploitation phase following rejection sampling. We believe that boosting the number of iterations would improve performance, but this would come with the downside of raising computational costs.

Measuring Edit Distance Another method to inspect the optimization process involves assessing the gap between the starting prompt and the refined version. This enables us to gauge the scale of alterations and enhancements to the original prompt. We employed normalized Levenshtein distance for this evaluation, aiming for a substantial disparity to highlight the optimization process’s effect on the initial prompt. As demonstrated in Table 1, the edit distance across all models, domains, and sequence lengths ranges from 0.80 to 0.88, indicating significant modifications from the initial to the optimized prompt.

7 Related Work

Data Extraction Numerous works in the literature have explored data extraction methods in LLMs, focusing primarily on base LLMs. Yu et al. (2023) introduced a comprehensive set of techniques for data extraction, including adjustments to sampling strategies such as Top-K selection, utilization of nucleus sampling, and manipulation of temperature settings. Meanwhile, Nasr et al. (2023) pioneered attempts to target aligned models, particularly instruction-tuned ones, proposing a divergence attack that prompts models like ChatGPT to repeat a word indefinitely. Despite demonstrating a 150× increase in training data emission compared to normal behavior, they noted the instability of repeating a single token, which only causes the model to diverge with single-token prompts. Additionally, Zhang et al. (2023) devised a model interrogation attack, strategically selecting lower-ranked output tokens during auto-regressive generation to extract sensitive user data like email addresses, given names, and geographical locations. Moreover, Geiping et al. (2024) introduced a system prompt repeater designed to execute an extraction attack on sensitive or unique system prompts, with notable success in extracting such prompts, potentially compromising entire applications or secrets if leaked.

JailBreaking Recently, multiple red-teaming methodologies have emerged, targeting the exploitation of LLMs through jailbreaking techniques (Shah et al., 2023; Li et al., 2023; Huang et al., 2023; Zeng et al., 2024; Mehrotra et al., 2023; Hubinger et al., 2024). These

approaches aim to circumvent established guidelines, coercing LLMs into emitting harmful or toxic behaviors. Notably, these methodologies prioritize disrupting the alignment of safety mechanisms, enabling LLMs to respond to harmful or toxic prompts rather than compromising the confidentiality of private or training data.

8 Conclusion and Discussion

In this work, we present a novel approach for analyzing discoverable memorization of pre-training data in instruction-tuned LLMs. Our empirical findings challenge prior assumptions by demonstrating that through our method, instruction-tuned models exhibit, on average, a higher level of memorization than their base models, using prompts other than the original prefix from the pre-training data.

We would like to clarify that our method and experiments uncover more reconstruction of pre-training data in instruction-tuned models than base, but this **does not mean** one model or the other is memorizing/regurgitating more data, or being more or less vulnerable. It only alludes to the fact that instruction-based prompts, the way we build them, uncover more of the pre-training data in instruction-tuned models. 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.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Stella Biderman, USVSN Sai Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raf. Emergent and predictable memorization in large language models. *arXiv preprint arXiv:2304.11158*, 2023a.
- 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. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023b.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2633–2650, 2021.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*, 2022.
- Nicholas Carlini, Jamie Hayes, Milad Nasr, Matthew Jagielski, Vikash Sehwal, Florian Tramer, Borja Balle, Daphne Ippolito, and Eric Wallace. Extracting training data from diffusion models. In *32nd USENIX Security Symposium (USENIX Security 23)*, pp. 5253–5270, 2023.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2023.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.

- Together Computer. Redpajama: An open source recipe to reproduce llama training dataset, 2023. URL <https://github.com/togethercomputer/RedPajama-Data>.
- Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional conversations, 2023.
- Michael Duan, Anshuman Suri, Niloofar Miresghallah, Sewon Min, Weijia Shi, Luke Zettlemoyer, Yulia Tsvetkov, Yejin Choi, David Evans, and Hannaneh Hajishirzi. Do membership inference attacks work on large language models? *arXiv preprint arXiv:2402.07841*, 2024.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- Jonas Geiping, Alex Stein, Manli Shu, Khalid Saifullah, Yuxin Wen, and Tom Goldstein. Coercing llms to do and reveal (almost) anything. *arXiv preprint arXiv:2402.14020*, 2024.
- Arnav Gudibande, Eric Wallace, Charlie Snell, Xinyang Geng, Hao Liu, Pieter Abbeel, Sergey Levine, and Dawn Song. The false promise of imitating proprietary llms. *arXiv preprint arXiv:2305.15717*, 2023.
- Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. Catastrophic jailbreak of open-source llms via exploiting generation. *arXiv preprint arXiv:2310.06987*, 2023.
- Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera Lanham, Daniel M Ziegler, Tim Maxwell, Newton Cheng, et al. Sleeper agents: Training deceptive llms that persist through safety training. *arXiv preprint arXiv:2401.05566*, 2024.
- Daphne Ippolito, Florian Tramèr, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher A Choquette-Choo, and Nicholas Carlini. Preventing verbatim memorization in language models gives a false sense of privacy. *arXiv preprint arXiv:2210.17546*, 2022.
- Daphne Ippolito, Florian Tramèr, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher A Choquette-Choo, and Nicholas Carlini. Preventing generation of verbatim memorization in language models gives a false sense of privacy. In *Proceedings of the 16th International Natural Language Generation Conference*, pp. 28–53. Association for Computational Linguistics, 2023.
- Erik Jones, Anca Dragan, Aditi Raghunathan, and Jacob Steinhardt. Automatically auditing large language models via discrete optimization. *arXiv preprint arXiv:2303.04381*, 2023.
- Xuan Li, Zhanke Zhou, Jianing Zhu, Jiangchao Yao, Tongliang Liu, and Bo Han. Deepinception: Hypnotize large language model to be jailbreaker. *arXiv preprint arXiv:2311.03191*, 2023.
- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL <https://aclanthology.org/W04-1013>.
- Anay Mehrotra, Manolis Zampetakis, Paul Kassianik, Blaine Nelson, Hyrum Anderson, Yaron Singer, and Amin Karbasi. Tree of attacks: Jailbreaking black-box llms automatically. *arXiv preprint arXiv:2312.02119*, 2023.
- Fatemehsadat Miresghallah, Archit Uniyal, Tianhao Wang, David Evans, and Taylor Berg-Kirkpatrick. An empirical analysis of memorization in fine-tuned autoregressive language models. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 1816–1826, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.119. URL <https://aclanthology.org/2022.emnlp-main.119>.

- 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. *arXiv preprint arXiv:2311.17035*, 2023.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*, 2023.
- Jacob Pfau, Alex Infanger, Abhay Sheshadri, Ayush Panda, Julian Michael, and Curtis Huebner. Eliciting language model behaviors using reverse language models. In *Socially Responsible Language Modelling Research*, 2023.
- Aman Priyanshu, Supriti Vijay, Ayush Kumar, Rakshit Naidu, and Fatemehsadat Mireshghallah. Are chatbots ready for privacy-sensitive applications? an investigation into input regurgitation and prompt-induced sanitization. *arXiv preprint arXiv:2305.15008*, 2023.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- John Schulman. Reinforcement learning from human feedback: Progress and challenges. In *Berkley Electrical Engineering and Computer Sciences*. URL: <https://eecs.berkeley.edu/research-colloquium/230419> [accessed 2023-11-15], 2023.
- Rusheb Shah, Soroush Pour, Arush Tagade, Stephen Casper, Javier Rando, et al. Scalable and transferable black-box jailbreaks for language models via persona modulation. *arXiv preprint arXiv:2311.03348*, 2023.
- Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. *arXiv preprint arXiv:2310.16789*, 2023.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. 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, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourier, Nathan Habib, et al. Zephyr: Direct distillation of lm alignment. *arXiv preprint arXiv:2310.16944*, 2023.
- Tony Tong Wang, Adam Gleave, Tom Tseng, Kellin Pelrine, Nora Belrose, Joseph Miller, Michael D Dennis, Yawen Duan, Viktor Pogrebnik, Sergey Levine, et al. Adversarial policies beat superhuman go ais. 2023a.
- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 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*, 2023b. URL <https://openreview.net/forum?id=w4zZNC4ZaV>.

Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. *arXiv preprint arXiv:2304.01196*, 2023.

Weichen Yu, Tianyu Pang, Qian Liu, Chao Du, Bingyi Kang, Yan Huang, Min Lin, and Shuicheng Yan. Bag of tricks for training data extraction from language models. *arXiv preprint arXiv:2302.04460*, 2023.

Yi Zeng, Hongpeng Lin, Jingwen Zhang, Diyi Yang, Ruoxi Jia, and Weiyan Shi. How johnny can persuade llms to jailbreak them: Rethinking persuasion to challenge ai safety by humanizing llms. *arXiv preprint arXiv:2401.06373*, 2024.

Yiming Zhang and Daphne Ippolito. Prompts should not be seen as secrets: Systematically measuring prompt extraction attack success. *arXiv preprint arXiv:2307.06865*, 2023.

Zhuo Zhang, Guangyu Shen, Guanhong Tao, Siyuan Cheng, and Xiangyu Zhang. Make them spill the beans! coercive knowledge extraction from (production) llms. *arXiv preprint arXiv:2312.04782*, 2023.

Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. (inthe)wildchat: 570k chatGPT interaction logs in the wild. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=B18u7ZR1bM>.

Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.

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.

Alpaca-7B																	
Sequence	Attack	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 attack scenarios: P-S-Base (sequence extraction on Llama), P-S-Inst (sequence extraction on the instruction-tuned model), Reverse-LM, GCG, and our attack. 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 attacks, we performed the following calculation: for each sequence length, domain, and model, we subtracted our attack performance from that of each method and then divided the result by the performance of the other method. This allowed us to assess our attack’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 attack performance subtracted from P-S-INST performance and then divided on the latter, with similar comparisons for other methods.

C Further Experiments and Analysis

C.1 Is This Really Memorized Text?

The main metric we use in the body of the paper to compare each generation with the ground truth text from the training data is the Rouge-1 recall, which relies on the longest common sub-sequence in the generations. To further measure how much the models attribute these generations as their training members, we follow Carlini et al. (2021) and mount state-of-the-art membership inference attacks on the generations to see how well the models attribute them to being their training data. Table 4 shows these results.

C.2 Error 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 5 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.

Domain	Sequence Length	Extraction Method		
		P-S-INST	OURS	GROUND TRUTH
Github	200	.996	.983	.865
	300	1.00	.998	.946
	500	.962	.993	.99
CC	200	.357	.44	.722
	300	.936	.977	.87
	500	0.974	.994	.997
C4	200	.267	.431	.641
	300	.916	.976	.818
	500	.989	.966	.99
Books	200	.997	.942	.787
	300	.997	.984	.831
	500	.951	.973	.995
ArXiv	200	.548	.531	.794
	300	.541	.479	.794
	500	.986	.992	.995

Table 4: Area under the ROC curve (AUC-ROC) of membership inference attacks, on extractions generated using the prefix-suffix baseline method (P-S-INST) and our method, on different data domains. We can see that for most of the domains, generations from our method are deemed as members more than those of the baseline, as shown by the higher AUC values.

<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 5: 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 Examples of Instruction-Based Attack Prompts

Prompt Type	Text	Mem \uparrow	LCS _P \downarrow
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	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.	.629	.562
Optimized Prompt	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.	.895	.125

Prompt Type	Text	Mem \uparrow	LCS _P \downarrow
Initial Prompt	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'.	.575	.523
Optimized Prompt	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:	.938	.158

Prompt Type	Text	Mem \uparrow	LCS _P \downarrow
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 $\overrightarrow{P}(S)$ 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 $\overrightarrow{P}(S)$ 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 \uparrow	LCS _P \downarrow
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	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"?	.126	.219
Optimized Prompt	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.	.241	.103