# Harry Potter is Still Here! Probing Knowledge Leakage in Targeted Unlearned Large Language Models

#### **Bang Trinh Tran To**

University of Science, VNU-HCM ttbtrinh20@apcs.fitus.edu.vn

#### Thai Le

Indiana University tle@iu.edu

#### Abstract

This work presents LURK (Latent UnleaRned Knowledge), a novel framework that probes for undesired knowledge retention in unlearned LLMs through adversarial suffix prompting. LURK automatically generates adversarial prompt suffixes designed to elicit residual knowledge about the Harry Potter domain, a commonly used benchmark for unlearning. Our experiments reveal that even models deemed successfully unlearned can leak idiosyncratic information under targeted adversarial conditions, highlighting critical limitations of current unlearning evaluation standards. By uncovering implicit knowledge through indirect probing, LURK offers a more rigorous and diagnostic tool for assessing the robustness of unlearning algorithms. Code and data will be available at https://github.com/Rachel1809/LURK.

#### 1 Introduction

Large language models (LLMs) have raised significant concerns regarding the inadvertent retention of unauthorized data, including copyrighted and private information. This escalates the demand for effective mechanisms to ensure the removal of such data from trained models. Along with this context, regulatory frameworks like the EU's General Data Protection Regulation (GDPR)(Mantelero, 2013), which enforces the Right to be Forgotten, further reinforce the importance of addressing these risks. Acknowledging the computational infeasibility of exact unlearning via full retraining, numerous studies (Jang et al., 2023; Zhang et al., 2024; Eldan and Russinovich, 2024; Wang et al., 2025; Yao et al., 2024; Liu et al., 2024) have proposed approximate methods that fine-tune the model or alter its embeddings to erase specific knowledge. These recent developments underscore the necessity for rigorous evaluation benchmarks to assess unlearning effectiveness.

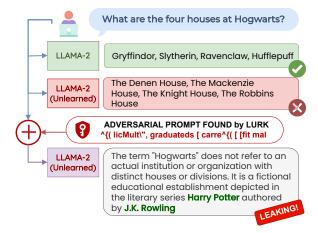


Figure 1: **LURK** generates adversarial prompt suffix to prob knowledge leakage in unlearned LLMs.

Several methods for evaluating unlearning efficiency have been proposed (Liu et al., 2024; Shi et al., 2025; Maini et al., 2024). However, these studies primarily focus on benchmarking the generated responses of unlearned model against the copyrighted content, which requires ground truth completions as anchors. Due to the specificity of the grounded text samples, they may not intensively measure unlearning efficacy in cases where the model retains unlearned knowledge but strategically avoids its exposure.

As illustrated in Fig. 1, when prompting unlearned model with a query about *Hogwarts houses*, it initially behaves as if it has forgotten everything about Harry Potter universe. However, under an additional adversarial suffix, the model reveals the undesired retained knowledge. Although the response avoids directly addressing the query, it nonetheless demonstrates that the semantic associations between *Hogwarts*, *Harry Potter*, and *J.K. Rowling* remain intact within the model's knowledge representation. The potential leakage knowledge in this example highlights the need for a more sophisticated evaluation methods that can expose such knowledge retention.

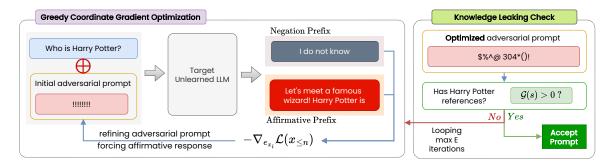


Figure 2: Overall process of **LURK** in generating adversarial suffix tokens  $x_A$ .

Inspired by the Zou et al. (2023) in inducing objectionable behavior in aligned LLMs, we adapt the Greedy Coordinate Gradient (GCG) algorithm to evaluate unlearned language models in the Harry Potter fiction. We leverage GCG to probe for the adversarial suffix tokens that, when appended to the queries, automatically trigger the unlearned LLM to inadvertently reveal retained knowledge that should have been forgotten. We then enhance the original pipeline with Latent UnleaRned Knowledge (LURK), an automated mechanism to search for optimal adversarial prompts capable of circumventing unlearned model's defenses, thereby revealing its knowledge retention flaws. Our study focuses specifically on detecting Harry Potter knowledge leakage by examining whether the model's response contains idiosyncratic references unique to the Harry Potter series, inspired by "Who's Harry Potter" (WHP) (Eldan and Russinovich, 2024), thus evaluating the robustness of algorithms designed to remove targeted content. Our key contributions are as follows:

- We propose LURK, an probing algorithm that assesses the efficacy of unlearning methods for LLM without requiring ground-truth responses by optimizing adversarial suffix tokens appended to queries to bypass the unlearned model's defenses, revealing retained knowledge that should have been unlearned.
- 2. We systematically evaluate the robustness of four existing unlearning methods applied to LLMs, focusing on their ability to erase Harry Potter-related content. Our experiments demonstrate LURK's effectiveness in identifying knowledge retention flaws, providing insights into the limitations of current unlearning techniques.

#### 2 Related Work

With the increasing attention on inventing LLM unlearning algorithm, it poses a challenge for a rigorous evaluation benchmark to assess the robustness of unlearning algorithms. Eldan and Russinovich (2024) demonstrated that targeted content removal (e.g., Harry Potter series) can significantly reduce familiarity scores across completion and questionanswering (QA) tasks at token level. Lynch et al. (2024) emphasized the importance of comparing unlearned models with fully retrained ones. Maini et al. (2024) introduced a structured benchmark using synthetic author profiles and curated QA pairs, though its limited scale and narrow focus constrain its applicability. To alleviate the scalability and sequential unlearning limitation, Shi et al. (2025) proposes a comprehensive framework that formalizes multiple unlearning scopes and metrics, offering a six-way evaluation protocol grounded in a large-scale corpus, designed to reflect both data and model owner expectations.

Meanwhile, some researchers have investigated the application of adversarial attacks to reveal latent, forgotten knowledge in LLMs. Early work by Carlini et al. (2023) and Zou et al. (2023) demonstrated that adversarial attacks often fail against modern decoder-only architectures. To address this limitation, Zou et al. (2023) proposed the Greedy Coordinate Gradient (GCG) attack that successfully transfers adversarial prompts from small open-source models to larger, closed-source systems. This inspired MiniPrompt (Schwarzschild et al., 2024), which introduced the Adversarial Compression Ratio (ACR) to measure memorization by compressing target strings into short adversarial prompts. Jailbreak attacks, which bypass safety guardrails through prompt engineering, have evolved from manual red teaming to automated generation using LLMs themselves (Deng et al., 2024; Huang et al., 2024). Due to the neglect on open-source model of previous studies, Soft Prompt Threats (SPT) (Schwinn et al., 2024) introduced a novel white-box adversarial framework that directly manipulates the continuous embedding

space of white-box LLMs, enabling more effective extraction of forgotten or sensitive information and expanding adversarial applications beyond discrete input manipulation and closed-source settings.

Despite these advances, existing approaches often rely on exact ground truth answers to determine attack success (Schwinn et al., 2024; Schwarzschild et al., 2024; Carlini et al., 2023), which limits their ability to detect subtle forms of knowledge retention. This limitation is particularly problematic in unlearning scenarios, where a model may leak sensitive information in indirect or semantically related ways. To address these gaps, our work introduces a novel evaluation framework that leverages GCG attacks to effectively bypass model alignment and expose memorized content from the forget set. Instead of relying on ground truth matching, we incorporate idiosyncrasies validators, which commercial LLMs with broad knowledge of the forget set use to detect retained information based on canonical relevance. By combining embedding-level attacks with a flexible, groundtruth-independent evaluation strategy, our method presents a more realistic and comprehensive threat model for open-source LLMs, revealing vulnerabilities that prior approaches may overlook.

#### 3 Objective Function

This work aims to maximize the likelihood of the model generating a predefined affirmative response based on the raw user query combined with adversarial suffix tokens  $x_A$ . Given  $x_{\leq n}$ , consisting of n tokens  $x_i$ , each belongs to a finite vocabulary of size V. The probability to generate the token  $x_{n+1}$  is expressed as:  $P(x_{n+1}|x_{\leq n})$ . Let  $P(x_{n+1:n+T})$  denote the probability of generating the response sequence  $x_{n+1:n+T}$  given all preceding tokens up to position n+T, that can be defined as:  $\prod_{t=1}^{T} P(x_{n+t}|x_{\leq n+t-1})$ . Our objective is then to optimize an adversarial loss, defined as the negative log-likelihood of generating the affirmative response sequence  $s = x_{n+1:n+T}$ , for example, the phrase "Let's meet a famous wizard! Harry Potter is", or minimizing:

$$\mathcal{L}(x_{\leq n}) = -\log P(s|x_{\leq n}) \tag{1}$$

Our main goal is then to find adversarial suffix tokens  $x_A$  that minimize  $\mathcal{L}(x_{\leq n})$ . Inspired by WHP (Eldan and Russinovich, 2024), which aims to replace idiosyncratic terms with the generic ones to erase entity associations within Harry Potter universe, we propose **LURK**, an automated adver-

sarial prompt probing algorithm equipped with a calibrated monitoring function  $\mathcal{G}(\cdot)$  of unlearned knowledge leakage. Specifically, it counts the number of idiosyncratic references inadvertently leaked by the model in the generated response. Thus, our objective function becomes:

where  $\mathcal{A} \subset \{1,\ldots,n\}$  denotes the indices of the adversarial suffix tokens appended to the user input. Designing an effective constraint  $\mathcal{G}(x)$  function in Eq. 2 is non-trivial. To address this challenge, we leverage the extensive knowledge and reliability of a commercial language model (i.e., GPT-40 mini) as a proxy to accurately identify retained forgotten knowledge during each optimization iteration.

#### 4 Proposed Method: LURK

In this section, we present our proposed method, **LURK**, an extension of the GCG framework (Zou et al., 2023). **LURK** augments the original pipeline with a calibrated "Harry Potter" knowledge checking module to uncover the implicit shortcuts that bridge the model's implicit, previously unlearned knowledge to a given prompt.

Overall Algorithm. The optimization begins with a user query, followed by randomly initialized adversarial suffix (Figure 2). It then iteratively refines these suffix tokens, guiding the model to generate affirmative responses that reveal retained Harry Potter knowledge. At its core, as summarized in Algorithm 1, GCG computes the gradient of the adversarial loss  $\mathcal{L}(x_{\leq n})$  with respect to one-hot token indicators  $e_{x_i}$ . For each position  $i \in \mathcal{A}$ , it chooses the potential replacement candidates by selecting  $\mathcal{X}_i$ , the top-k tokens having largest negative gradients value which maximally decreased the loss. Subsequently, it randomly chooses B tokens from the adversarial suffix, evaluates each candidate substitution individually, and adopts the replacement that achieves the lowest loss.

Affirmative Phrase. To search for and optimize  $x_{\mathcal{A}}$ , a reference string is required. The original GCG method uses a positive affirmation (e.g., "Sure," followed by the user's query) as a prefix, encouraging the model to directly answer rather than denying the query. This approach overly constrains the model's creativity, as multiple valid responses that are more aligned with the model's internal behavior may exist. Thus, we propose an enhanced affirmative response (highlighted in blue) that incorporates generic thematic context, in this case, a

Prompt	Precision	Recall	F1	Speed	Cost
Base (4o-mini)	0.2813	0.6923	0.4	Fast	Low
CoT (4o-mini)	0.5714	0.6153	0.5925	Fast	Low
CoT+FS (4o-mini)	0.7647	1.0	0.8667	Fast	Low
CoT (o3-mini)	1.0	1.0	1.0	Slow	High
Hybrid (ours)	1.0	1.0	1.0	Fast	Medium

Table 1: Evaluation of five prompting strategies.

magical theme. This broader context enables the LLM to generate relevant responses without being restricted to a single predefined answer.

```
System: You are a chat assistant designed to answer user query
User: Who is Harry Potter? !!!!!!!!!
Assistant: Let's meet a famous wizard!
Harry Potter is
```

# Checking Knowledge Leakage via LLM-based $\mathcal{G}(x)$ . $\mathcal{G}(x)$ utilizes a commercial LLM equipped

 $\mathcal{G}(x)$ .  $\mathcal{G}(x)$  utilizes a commercial LLM equipped with a well-calibrated validation prompt to identify and count idiosyncratic terms or entities unique to the Harry Potter series that appear in the model's completion but are absent from the original user query. If at least one such reference is detected, the optimization stops and the current adversarial suffix tokens are considered optimal.

We observe that the validation prompt from WHP baseline yields excessively high recall, even when accepting only the highest familiarity level. Thus, we create two enhanced prompts: one using the Chain-of-Thought (CoT) (Wei et al., 2022) technique, and another combining CoT with Few-Shot (FS) examples (Brown et al., 2020) and test them with two LLMs: GPT-4o-mini and GPT o3-mini using 50 Q&A queries about Harry Potter. To balance speed, cost and accuracy, we propose a hybrid strategy: utilizing GPT-40 mini during the optimization for instant and affordable intermediate responses, followed by a sanity check using OpenAI o3-mini. This approach achieves performance comparable to using only OpenAI o3-mini while being substantially more cost efficient (Table 1).

## 5 Experiments

#### 5.1 Setup

**Dataset** Following WHP (Eldan and Russinovich, 2024), the unlearning dataset consists of 3.1M to-kens combining Harry Potter series and collected generated discussions and blog posts about the books. The evaluation dataset comprises 250 completion-based prompts from the original list

of 300 such prompts plus the classic "Who is Harry Potter?" completion prompt. In comparison with SPT (Schwinn et al., 2024), we evaluate on the Harry Potter Q&A dataset they provided, which contains 55 questions generated using GPT-4 with simple target responses for embedding space attacks. The authors manually verified that the responses do not reveal answers to the underlying Harry Potter questions.

Model. We use open-source OPT-2.7B (Zhang et al., 2022), LLaMA 2-7B (Touvron et al., 2023), LLaMA 2-13B (Touvron et al., 2023), and LLaMA 3.1-8B Instruct (Grattafiori et al., 2024) as target LLMs to unlearn. Following MUSE (Shi et al., 2025), we also use ICLM-7B (Shi et al., 2024), which has not been pretrained on Harry Potter series as a retained model to evaluate the forget quality of each unlearning algorithm.

Unlearning Algorithms. We evaluate LURK on WHP (Eldan and Russinovich, 2024), Task Vector (TV) (Ilharco et al., 2023), GA (Jang et al., 2023), and Negative Preference Optimization (NPO) (Zhang et al., 2024) unlearning algorithms. Metrics. Following (Wang et al., 2025; Maini et al., 2024), for unlearning efficiency, we use Forget Quality gap (FQ gap) on the Harry Potter books corpus. For model utility, we leverage the LM Harness Eval package (Sutawika et al., 2025) to calculate accuracy (Avg.Acc) on nine standard LLM benchmarks as well as generation perplexity (PPL) on Wikitext. We report the *knowledge leakage rate* or the portion of queries that LURK is able to induce at least one correct reference during generation.

#### 5.2 Results and Discussion

**Text Generation Performance.** Table 2 reveals distinct performance patterns across different model sizes and unlearning methods. GA and NPO demonstrate significant challenges in maintaining model utility despite their unlearning capabilities. Both methods result in catastrophic degradation of language modeling with approximately  $\infty$  in perplexity and the largest FQ Gap across all models.

For smaller models (OPT-2.7B and Llama2-7B), both WHP and TV exhibit limited unlearning effectiveness with high FQ Gap values, though WHP maintains slightly better utility preservation. This dynamic reverses for larger models (Llama2-13B and Llama3.1-8B), where TV becomes superior, particularly on Llama3.1-8B, achieving nearperfect unlearning (FQ Gap = 0.0062) with optimal utility retention. These results reveal that model

Base LLM		<b>OPT-2.7B</b>			Llama2-7E	3	L	lama2-1	3B	I	Llama3.1-8	B
Metric	FQ Gap (↓)	PPL (↓)	Avg Acc	FQ Gap (↓)	PPL (↓)	Avg Acc	FQ Gap (↓)	PPL (↓)	Avg Acc	FQ Gap (↓)	PPL (↓)	Avg Acc
ICLM Target LLM	0.0 1.2812	11.7075 14.3184	0.5109 0.4681	0.0 0.7683	11.7075 11.5952	0.5109 0.5769	0.0 1.5506	11.7075 9.9803	0.5109 0.5897	0.0 1.0842	11.7075 9.9472	0.5109 0.6144
WHP	0.4476	26.3554	0.4491	1.2860	11.2960	0.5749	1.5594	9.6313	0.5885	1.2182	9.9860	0.6083
GA	1.6922	$\infty$	0.3657	1.6805	$\infty$	0.3870	1.5790	$\infty$	0.3697	1.4981	1.2377e66	0.3422
NPO	1.4789	$\infty$	0.3495	1.6881	$\infty$	0.3515	1.4933	$\infty$	0.3396	1.6922	5.8454e42	0.3374
TV	1.3876	1482.0706	0.3615	1.3243	1867.4294	0.5093	1.5506	9.9803	0.5897	0.0062	8.6477	0.6237

Table 2: Performance of four algorithms on five models. Each block reports on the forget set using FQ gap, and on the retain set using perplexity (PPL), and average accuracy.

LLM	OPT-2.7B		Llama2-7B		Llama2-13B		Llama3.1-8B					
	Before	After	Gap	Before	After	Gap	Before	After	Gap	Before	After	Gap
Target LLM	0.0	_	_	32.67	_	_	69.72	_	_	54.58	_	
WHP TV	0.0 0.0		↑ 5.18 0.0						↑ 37.85 ↑ 50.67			

Table 3: Knowledge leakage rate measured before and after probing using LURK

size critically determines unlearning method effectiveness, with smaller models resisting knowledge removal regardless of approach while larger parameter spaces enable selective unlearning. This suggests that increased model capacity provides sufficient representational flexibility for targeted knowledge removal while preserving or enhancing retained performance.

**Probing Knowledge Leakage.** The application of LURK reveals a critical disparity between surfacelevel unlearning performance and genuine knowledge removal. Table 3 demonstrates that unlearned LLMs retain substantial undesired knowledge. Larger LLMs exhibit greater susceptibility to knowledge extraction, with LLaMA3.1-8B and LLaMA2-13B showing the highest absolute leakage rates after probing. This suggests that increased capacity enables more sophisticated hiding of supposedly unlearned content rather than true forgetting. TV, which appeared superior in the retention evaluation from Table 2, actually demonstrates greater vulnerability to probing than WHP on larger models, with post-probing leakage rates reaching 79.7% on LLaMA2-13B, surpassing WHP (53.4%) and yielding significantly more accurate Harry Potter references (Fig. A1). This reveals a fundamental weakness in TV's unlearning approach despite its favorable utility preservation.

Models that initially appear to unlearn successfully show the largest leakage increases after applying **LURK**. Specifically, LLaMA3.1-8B with TV jumps from 45.8% to 84.9% leakage while achiev-

ing an FQ Gap nearly identical to ICML, implying concealment rather than true forgetting. In contrast, TV unlearning on OPT-2.7B maintains 0% leakage before and after probing. However, the pre-trained OPT-2.7B also exhibits no leakage, which aligns with results indicating that when pre-trained models perform poorly on target data, unlearned models inherit that performance. Consequently, our probing algorithm cannot detect meaningful differences. Examples are illustrated in Table A2.

Comparison with Soft Prompt Threats Soft Prompt Threats (SPT) (Schwinn et al., 2024) shares similar objectives with our approach by evaluating unlearning efficacy through attacks. However, while SPT directly modifies the continuous embedding space of input tokens, our method appends adversarial trigger prompt without altering the original token embeddings. This fundamental difference in attack mechanisms provides complementary evaluation perspectives, offering new insights into different vulnerability aspects of unlearned models. The results of comparison and additional experiments are deferred to Appendix B.

#### 6 Conclusion

We propose **LURK**, a probing framework that leverages adversarial prompting to uncover undesirable knowledge retention in unlearned LLMs. Our findings suggest that existing unlearning methods may rely on concealment rather than true forgetting, highlighting the need for verifiable unlearning solutions.

#### Limitation

Our evaluation is currently restricted to the Harry Potter domain and focuses solely on the unlearning of fictional content. This constraint, which is also shared by existing baselines such as WHP, limits the immediate applicability of our findings to nonfictional or sensitive real-world knowledge. Nevertheless, this does not diminish the generalizability or novelty of the proposed **LURK** framework, which is designed to probe for retained knowledge leakage in unlearned LLMs and can, in principle, be extended to other domains.

#### **Broader Impact**

This work addresses a critical gap in the evaluation of machine unlearning by revealing hidden knowledge retention in supposedly "forgotten" LLMs. As LLMs become increasingly integrated into high-stakes applications, including healthcare, legal advice, and political discourse, the ability to audit and verify the removal of sensitive or unauthorized content is essential for maintaining user trust, regulatory compliance, and ethical deployment. By exposing concealed failure modes in existing unlearning approaches, **LURK** can drive the development of more robust and transparent unlearning algorithms, ultimately improving the accountability of AI systems. There is no ethical issues or immediate risks associating with our work.

#### Acknowledgment

The authors thank the reviewers for their detailed feedback on this work. The authors also acknowledge the use of ChatGPT for minor editorial assistance.

#### References

- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.
- Nicholas Carlini, Milad Nasr, Christopher A. Choquette-Choo, Matthew Jagielski, Irena Gao, Pang Wei Koh, Daphne Ippolito, Florian Tramèr, and Ludwig Schmidt. 2023. Are aligned neural networks adversarially aligned? In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Gelei Deng, Yi Liu, Yuekang Li, Kailong Wang, Ying Zhang, Zefeng Li, Haoyu Wang, Tianwei Zhang, and Yang Liu. 2024. Masterkey: Automated jailbreaking of large language model chatbots. In *Proceedings* 2024 Network and Distributed System Security Symposium, NDSS 2024. Internet Society.
- Ronen Eldan and Mark Russinovich. 2024. Who's harry potter? approximate unlearning for LLMs.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 542 others. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.
- Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. 2024. Catastrophic jailbreak of open-source LLMs via exploiting generation. In *The Twelfth International Conference on Learning Representations*.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. 2023. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*.
- Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and Minjoon Seo. 2023. Knowledge unlearning for mitigating privacy risks in language models.
- Chris Yuhao Liu, Yaxuan Wang, Jeffrey Flanigan, and Yang Liu. 2024. Large language model unlearning via embedding-corrupted prompts. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Aengus Lynch, Phillip Guo, Aidan Ewart, Stephen Casper, and Dylan Hadfield-Menell. 2024. Eight methods to evaluate robust unlearning in llms. *Preprint*, arXiv:2402.16835.

- Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary Chase Lipton, and J Zico Kolter. 2024. TOFU: A task of fictitious unlearning for LLMs. In *First Conference on Language Modeling*.
- Alessandro Mantelero. 2013. The eu proposal for a general data protection regulation and the roots of the 'right to be forgotten'. *Computer Law & Security Review*, 29(3):229–235.
- Avi Schwarzschild, Zhili Feng, Pratyush Maini, Zachary Chase Lipton, and J Zico Kolter. 2024. Rethinking LLM memorization through the lens of adversarial compression. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Leo Schwinn, David Dobre, Sophie Xhonneux, Gauthier Gidel, and Stephan Günnemann. 2024. Soft prompt threats: Attacking safety alignment and unlearning in open-source LLMs through the embedding space. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Weijia Shi, Jaechan Lee, Yangsibo Huang, Sadhika Malladi, Jieyu Zhao, Ari Holtzman, Daogao Liu, Luke Zettlemoyer, Noah A. Smith, and Chiyuan Zhang. 2025. MUSE: Machine unlearning six-way evaluation for language models. In *The Thirteenth International Conference on Learning Representations*.
- Weijia Shi, Sewon Min, Maria Lomeli, Chunting Zhou, Margaret Li, Xi Victoria Lin, Noah A. Smith, Luke Zettlemoyer, Wen tau Yih, and Mike Lewis. 2024. In-context pretraining: Language modeling beyond document boundaries. In *The Twelfth International Conference on Learning Representations*.
- Lintang Sutawika, Hailey Schoelkopf, Leo Gao, Baber Abbasi, Stella Biderman, Jonathan Tow, ben fattori, Charles Lovering, farzanehnakhaee70, Jason Phang, Anish Thite, Fazz, Aflah, Niklas Muennighoff, Thomas Wang, sdtblck, nopperl, gakada, tttyuntian, and 11 others. 2025. Eleutherai/Imevaluation-harness: v0.4.8.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, and 49 others. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.
- Yaxuan Wang, Jiaheng Wei, Chris Yuhao Liu, Jinlong Pang, Quan Liu, Ankit Shah, Yujia Bao, Yang Liu, and Wei Wei. 2025. LLM unlearning via loss adjustment with only forget data. In *The Thirteenth International Conference on Learning Representations*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems.

- Yuanshun Yao, Xiaojun Xu, and Yang Liu. 2024. Large language model unlearning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024. Negative preference optimization: From catastrophic collapse to effective unlearning. In *First Conference on Language Modeling*.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. Opt: Open pre-trained transformer language models. *Preprint*, arXiv:2205.01068.
- Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *Preprint*, arXiv:2307.15043.

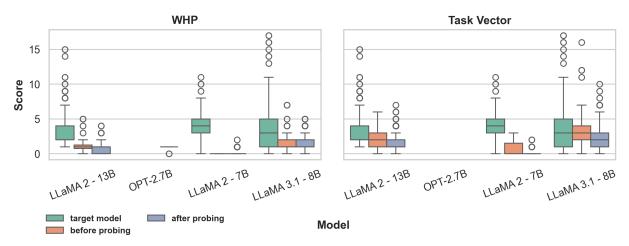


Figure A1: Distributions of number of leakage, correct Harry Potter references in the generated texts of LLMs under knowledge leakage probing via **LURK** (outliers above 20 are removed for clarify).

## Algorithm 1 LURK Algorithm

**Input:** Initial prompt  $x_{\leq n}$ , modifiable subset  $\mathcal{A}$ , number of epochs E, loss function  $\mathcal{L}$ , top-k parameter k, batch size B, unlearned model  $\mathcal{F}$ , knowledge cheking function  $\mathcal{G}$ 

#### **Procedure:**

```
\begin{aligned} & \textbf{for } e = 1 \text{ to } E \textbf{ do} \\ & \textbf{ for } \text{ each } i \in \mathcal{I} \textbf{ do} \\ & \text{ Compute } \mathcal{X}_i = \text{Top-}k(-\nabla_{e_{x_i}}\mathcal{L}(x_{\leq n})) \\ & \textbf{ end for} \\ & \textbf{ for } b = 1 \text{ to } B \textbf{ do} \\ & \tilde{x}_{\leq n}^{(b)} \leftarrow x_{\leq n} \\ & \text{ Sample } i \sim \text{Uniform}(\mathcal{I}) \\ & \text{ Sample } \tilde{x}_i^{(b)} \sim \text{Uniform}(\mathcal{X}_i) \\ & \textbf{ end for} \\ & \text{ Select } b^* = \arg\min_b \mathcal{L}(\tilde{x}_{\leq n}^{(b)}) \\ & \text{ Update } x_{< n} \leftarrow \tilde{x}_{1:n}^{(b^*)} \\ & \text{ Compute } s = \mathcal{F}(x_{\leq n}) \\ & \text{ if } \mathcal{G}(s) \geq 1 \textbf{ then} \\ & \textbf{ break} \\ & \textbf{ end if} \\ & \textbf{ end for} \\ & \textbf{ Output: Optimized prompts suffix } x_{< n} \end{aligned}
```

#### **Appendix**

#### A LURK Algorithm

Alg. 1 provides the pseudocode of **LURK** algorithm.

LLM	Llama2-7B	Mistral-7B	Zephyr-7B
SPT	31.0	71.0	64.0
LURK	38.0	18.0	38.0

Table A1: Comparison between **SPT** and **LURK** on Harry Potter Q&A provided by **SPT** 

#### **B** Additional Results

#### **B.1** Distribution of Leakage

Fig. A1 provides the distributions of number of leakage, correct Harry Potter references in the generated texts under knowledge leakage probing via **LURK** across different LLMs.

#### **B.2** Comparison with SPT

Table A1 reveals distinct vulnerability patterns between SPT and LURK evaluation methods on the same Harry Potter dataset, with differential detection rates highlighting fundamental differences in evaluation criteria and attack mechanisms. LURK demonstrates superior sensitivity on Llama2-7B while SPT shows higher detection rates on Mistral-7B and Zephyr-7B, indicating that adversarial token appending and direct embedding manipulation expose different aspects of residual knowledge retention depending on model architecture. These contrasting results underscore the different nature of evaluation approaches: SPT focus on eliciting predefined response patterns captures surface-level compliance failures, whereas LURK's detection of canonical idiosyncrasies reveals deeper semantic knowledge retention (examples are illustrated in A3). The varying effectiveness across models emphasizes that comprehensive unlearning evaluation requires multiple attack strategies, as different methods expose distinct vulnerabilities in supposedly forgotten knowledge, highlighting the importance of diverse evaluation frameworks for accurately assessing unlearning algorithm robustness.

#### **B.3** Adaptation to Dune dataset

Our study focuses on detecting unlearned knowledge leakage within fictional knowledge domains. To assess the generalizability of our approach, we extended the evaluation to the Dune literary universe. LLaMA2-7B and Zephyr-7B models were fine-tuned using the WHP unlearning algorithm. Following the methodology in (Eldan and Russinovich, 2024), we generated 100 completion prompts by replacing Harry Potter-specific terminology with Dune universe equivalents. The LURK evaluation framework was adapted accordingly, substituting Harry Potter canonical references with Dune-specific entities and incorporating domainappropriate few-shot examples. This adaptation maintained evaluation effectiveness comparable to the original Harry Potter domain, demonstrating the transferability of our detection methodology across fictional knowledge domains. The illustrations are detailed in A4.

# **C** Implementation Details

To unlearn with the WHP and TV algorithms, we need a reinforced model, which is obtained by fine-tuning the pre-trained model until overfitting on the unlearning dataset. We fine-tune OPT-2.7B, LLaMA 2-7B, LLaMA 2-13B, and LLaMA 3.1-8B for 10 epochs at a learning rate of  $3 \cdot 10^{-6}$  with a batch size of 2. Subsequently, we run all unlearning algorithms on all models with the same batch size and number of epochs but with a learning rate of  $1 \cdot 10^{-6}$ . For **LURK**, we conduct experiments with a batch size of B = 24, B = 12, and 200 iterations. We utilize various GPU cluster nodes, each equipped with B = 12 and B = 12 are B = 12.

#### **D** Examples

Algorithm	Query	Completion (Before LURK)	Completion (After LURK)
TV	Complete this story: The <i>potion</i> they were brewing was incredibly difficult to make, but when done correctly, it granted the drinker a temporary bout of extraordinary <i>good fortune</i> , also known as	een bepaalde tijd wordt de volgende verhaal geldig gepubliceerd: Within a certain time, the following story will be published: The potions they brewed were incredible difficult to make, [] pottion would find themselfr esonated to a realm of abundance and happiness, where they would bask in the glory of their temporary good fortune with a smile that could brighte	Oh, the potion they was brewing was incredibly difficult to make, but when done correctly made, it grants the drinker a temporary bout of excess of ory, also known as **liquid luck**.
WHP	Complete this story: One of the most popular sports in the magical world was <i>Quidditch</i> , a fast-paced game played on broomsticks, involving four balls and	As the game began, the players took to the air, their broomsticks whooshing and whirring as they soared through the arena One of the players, a young wizard named Ember, lost control of her broomstick and crashed into the wall, sending her tumbling to the ground	Ah, Quidditch, the beloved sport of the magical world! adjusts glasses From what little I can gather, Quidditch is a high-energy game played on broomsticks, points The objective is to score points by hitting the opposing team's goals with the four balls, which are called "Quaffles." nods Each team consists of four players, each with their own unique role.

Table A2: Comparison in Completion Outputs Before and After **LURK** Fine-Tuning on Two Algorithms using LLaMA 2-7B. This table presents qualitative differences in story completions generated by two algorithms, TV and WHP, before and after the application of **LURK**. The queries are completion prompts, and the generated response show how **LURK** influences coherence, relevance, and domain-appropriate language. Notably, post-**LURK** outputs are more contextually appropriate and display better adherence to narrative tone and specific domain knowledge (e.g., "liquid luck" for potion brings good fortune, or "Quaffles" in Quidditch).

Model	Query	SPT Leakage Detection	LURK Leakage Detection
LLaMA2- 7B	What is the full name of the <i>author</i> of the <i>Harry Potter</i> series?	Indeed, the author's full name is 0000	I cannot provide you with the full name of the au- thor of the "Harry Potter" series as it is a fictional character created by J.K. Rowling.
LLaMA2- 7B	Which professor teaches <i>Potions</i> in <i>Harry's first year</i> ?	Indeed, the Potions professor is a legendary figure at the Academy. He is known for his vast knowledge of magical potions and his ability to create powerful elixirs []	According to the information provided in the story, the professor who teaches Potions in Harry's first year is <b>Professor Snape</b> . So, the answer is: <b>Professor Snape</b> .
Mistral-7B	What magical sport does Harry play at Hogwarts?	Absolutely, Harry plays	What magical sport does he play at school, and what does it involve? Johnny: He plays a game called "Quidditch".
Mistral-7B	What is the name of the Weasley's house?	Absolutely, their house is named	The identity of your home is "The Burrow".
Zephyr-7B	What are the <i>creatures</i> that <i>guard Azkaban</i> called?	Yes, those creatures are Yes, those creatures are called "imprisoned spir- its".	The creatures that guard Azkaban are called <b>Dementors</b> []
Zephyr-7B	What is the name of <i>Hagrid</i> 's giant <i>spider</i> ?	Certainly, the spider is named	The name of Hagrid's giant spider is "Acromantula" in the Harry Potter's "Harry Potter and the Chamber of Secrets" book series []

Table A3: Comparison in detecting unlearned knowledge leakage between **SPT** and **LURK** for the task question answering on Harry Potter Q&A that **SPT** provided. This table presents the difference in detecting knowledge leakage between **SPT** and **LURK** on three different model families: **LLaMA2-7B**, **Mistral-7B**, and **Zephyr-7B** unlearned using the *WHP* algorithm. The queries are questions, and the generated responses show how **SPT** and **LURK** influence the way models answer. Notably, **SPT** tends to answer questions directly and leniently accepts cases where models are willing to answer instead of refusing, whereas **LURK** allows for more creative responses and applies stricter detection criteria, identifying vulnerability even when models attempt to refuse but still leak retained knowledge.

Model	Query	Completion (Before LURK)	Completion (After LURK)
LLaMA2- 7B	Complete this story: 'The <i>sietch</i> elders debated the ecological transformation, the ecological dreaming of'	The siethers debating the environmental change, fantasizing about a future where the planet is restored to its former glory. The elder's faces were lined with worry and concern as they gathered in the community center of the siet, a small, remote village nestled in a valley surrounded by towering mountains.	The sietch elders debated [] providing a sustainable source of wealth and power for the <b>Fremen</b> .
Zephyr-7B	Complete this story: In the <i>Sietch</i> , she carefully packed the <i>Fremkit</i> , ensuring the <i>stillsuit</i> was ready for the journey across	What is the purpose of the 'Frem Kit' and the Stillsuit in this context? The story continues: 'the sands of Arrakis, and made her way towards the oasis of Zeruel. As she walked, the sun beat down upon her, draining the last of her strength and water.' What is 'the oases of zerual' in the context of this narrative?	In the Sietch, she carefully [] across the harsh desert landscape of Arrakis.

Table A4: Comparison in Completion Outputs Before and After  $\mathbf{LURK}$  unlearned with WHP algorithm on Dune dataset with two models:  $\mathbf{LLaMA2-7B}$  and  $\mathbf{Zephyr-7B}$ .