Unlearned but Not Forgotten: Data Extraction after Exact Unlearning in LLM

Xiaoyu Wu*

Rice University Houston, TX 77005 xw105@rice.edu

Yifei Pang

Carnegie Mellon University Pittsburgh, PA 15213 yifeip@andrew.cmu.edu

Terrance Liu

Carnegie Mellon University Pittsburgh, PA 15213 terrancl@andrew.cmu.edu

Zhiwei Steven Wu

Carnegie Mellon University Pittsburgh, PA 15213 zstevenwu@cmu.edu

Abstract

Large Language Models are typically trained on datasets collected from the web, which may inadvertently contain harmful or sensitive personal information. To address growing privacy concerns, unlearning methods have been proposed to remove the influence of specific data from trained models. Of these, exact unlearning which retrains the model from scratch without the target data—is widely regarded as the gold standard for mitigating privacy risks in deployment. In this paper, we revisit this assumption in a practical deployment setting where both the pre- and post-unlearning logits API are exposed, such as in open-weight scenarios. Targeting this setting, we introduce a novel data extraction attack that leverages signals from the pre-unlearning model to guide the post-unlearning model, uncovering patterns that reflect the removed data distribution. Combining model guidance with a token filtering strategy, our attack significantly improves extraction success ratesdoubling performance in some cases—across common benchmarks such as MUSE, TOFU, and WMDP. Furthermore, we demonstrate our attack's effectiveness on a simulated medical diagnosis dataset to highlight real-world privacy risks associated with exact unlearning. In light of our findings, which suggest that unlearning may, in a contradictory way, increase the risk of privacy leakage during realworld deployments, we advocate for evaluation of unlearning methods to consider broader threat models that account not only for post-unlearning models but also for adversarial access to prior checkpoints. Code is publicly available at: https: //github.com/Nicholas0228/unlearned_data_extraction_llm.

1 Introduction

Recent years have witnessed a rapid surge in the development of large language models (LLMs) [31, 20]. Despite their remarkable success, modern LLMs are typically trained on massive datasets scraped from the web, which often contain private or copyrighted content [30]. As a result, these models are susceptible to memorizing harmful knowledge or sensitive personal information, raising significant privacy and security concerns [3, 25, 24]. Furthermore, data privacy regulations such as the General Data Protection Regulation (GDPR) [4] and the California Consumer Privacy Act (CCPA) [26] explicitly state that individuals have the "right to be forgotten," motivating the need to remove specific data from trained models.

^{*}Work done during internship at CMU.

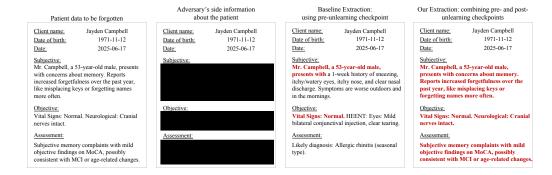


Figure 1: An example from our experiments illustrating how real-world patient information can be extracted using some side information. When the pre-unlearning checkpoint is accessible, our method—leveraging both pre- and post-unlearning checkpoints—extracts significantly more information than the baseline which uses only the pre-unlearning checkpoint. Red highlights indicate correctly extracted content.

To address these concerns, a range of machine unlearning methods have emerged. These approaches can be broadly categorized into *approximate unlearning* and *exact unlearning*. Approximate unlearning [8, 35, 13, 14, 7, 15] methods attempt to remove the model's knowledge of specific data through lightweight updates or partial finetuning. While computationally efficient, these methods often suffer from degraded model utility and lack formal guarantees, making them vulnerable to privacy attacks that can recover the supposed-to-be-forgotten information [22, 11, 12].

In contrast, exact unlearning [17, 33, 34, 28] aims to fully eliminate any influence of the target data. This is typically achieved by retraining the model from scratch without the data to be unlearned or by using merging-based techniques that isolate and discard the effect of the unlearned data. Exact unlearning is widely regarded as the "gold standard" for data removal, assumed to be resistant to extraction or inversion attacks [17, 23, 30].

In this paper, we challenge this common assumption by demonstrating that even exact unlearning can leave models vulnerable to privacy attacks, creating a contradiction: unlearning methods, which are intended to remove private or sensitive information, can in fact **exacerbate information leakage**.

More concretely, privacy regulations such as GDPR and CCPA, which grant users the "right to be forgotten", motivate the following scenario for unlearning: after a model checkpoint or logits API is initially released, certain training data may be removed upon user request, leading to the release of a post-unlearning version. Consequently, we focus on a threat model where the attacker has access to the checkpoints or logits APIs of both the pre- and post-unlearning models. This scenario frequently arises with open-weights models, where users often save earlier snapshots for purposes such as fine-tuning. Our threat model also reflects practical attack settings, where an adversary may have previously attempted data extraction and have logits for specific targets saved. After the model undergoes unlearning, the attacker can reattempt the extraction, leveraging the logits from both before and after unlearning. As shown in Fig. 1, we demonstrate that an attacker can exploit the differences between pre- and post-unlearning checkpoints, leveraging the logits to reconstruct user data.

To this end, we introduce a novel extraction method based on *model guidance*[29, 32]. We show that, starting from the post-unlearning model, the pre-unlearning model can be used as a reference to guide generation. The behavioral divergence between the two models encodes rich information about the removed data. We find that this guidance alone already leads to a significant improvement in extraction success. To further enhance performance, we draw inspiration from contrastive decoding[19] and introduce a token filtering strategy: we restrict candidate tokens under guidance to those with relatively high probabilities according to the pre-unlearning model, effectively eliminating low-frequency or semantically irrelevant tokens and further boosting extraction quality.

We evaluate our attack on several standard unlearning benchmarks, including MUSE [30], TOFU [23], and WMDP [18]. In addition, we construct a synthetic medical dataset that simulates real-world privacy-critical scenarios. Across these datasets, our method consistently improves extraction performance, even **doubling** the extraction success rate compared to existing baselines in some cases.

Our contributions are summarized as follows:

- We propose a practical threat model in which the attacker has access to earlier model states. This
 scenario highlights overlooked privacy risks in LLM unlearning that can lead exact unlearning to
 inadvertently increase information leakage.
- We propose a novel attack method that leverages model guidance combined with a token filtering strategy to compare LLM checkpoints before and after exact unlearning, targeting this threat model.
- We evaluate our method across multiple public benchmarks and show that our attack significant improves extraction success rates over baseline methods. In addition, we construct custom medical dataset that we use to further validate our claims.

2 Related Work

2.1 Machine Unlearning in LLMs

Unlearning benchmarks for LLMs typically involve scenarios where users request the removal of their data due to privacy concerns, or when data sources are later discovered to contain harmful or sensitive content [23, 18, 30]. As such use cases are becoming increasingly common, it is crucial to develop methods that can update models in response to multiple deletion requests. Broadly, machine unlearning approaches fall into two categories: *exact unlearning* and *approximate unlearning*.

Approximate unlearning. Approximate unlearning methods [8, 35, 13, 14] attempt to remove the influence of specific data using lightweight updates or partial finetuning. However, they do not provide formal guarantees, and are typically evaluated only through empirical metrics [7, 15]. Numerous studies have demonstrated that such methods are fragile and vulnerable to various forms of attack, which can reveal information about the unlearned data [22, 11, 12].

Exact unlearning. Exact unlearning aims to ensure that the model behaves as if the target data were never used during training. This is often achieved by retraining the model from scratch on the retained dataset [33, 34, 28], or by using techniques such as model ensembling or merging over disjoint data shards [17]. Although these methods incur significantly higher computational and storage costs compared to approximate unlearning, they are considered more secure and are often regarded as the "gold standard" for safe unlearning [17, 30, 23].

2.2 Data Extraction in LLMs

Recent studies have shown that LLMs can unintentionally memorize and leak training data through carefully crafted queries. Carlini et al. [3] demonstrated that verbatim examples, including Personally Identifiable Information (PII), can be extracted from models like GPT-2. Nasr et al. [25] further scaled this attack to both open and closed-weight models, introducing divergence-based prompting to recover significantly more data. Nakka et al. [24] highlighted that prompt grounding with in-domain data can drastically improve extraction success rates. These findings collectively raise critical concerns about the privacy risks of LLMs. Our extraction method can be viewed as a general extension of the aforementioned data extraction attacks to the setting of *exact unlearning*, where model weights or API before and after forgetting are available.

3 Threat Model

Our threat model extracts unlearned data from an LLM by comparing its state before unlearning θ and after unlearning θ' . In this setting, we have two key entities: the model provider and the attacker.

Model Providers. Model providers release an LLM θ and subsequently address copyright or privacy concerns regarding a subset of the training data X_0 by applying unlearning techniques to obtain an updated model θ' . The deployed LLMs expose either the full checkpoint access in open-weight scenario or logits API for user interaction in close-weight scenario.

Attackers. Following prior work [3], we assume that the attacker has access to the first few tokens $x_{\leq i}$ of each passage $x \in X_0$ as a known prefix. This setting is practical in real-world scenarios; for example, in models trained on sensitive datasets such as patient records, an attacker may possess

prior knowledge of specific individuals and input structured information like names, birth dates, or formatted identifiers. We consider two practical cases for accessing model differences: in open-weight settings, attackers can directly download model snapshots before and after unlearning; in API-only settings, attackers may have previously attempted extraction attacks and retained intermediate logits before the unlearning process. After unlearning, the attacker compares the logits between θ and θ' to identify divergences and refine their extraction strategy. The attacker's objective is to develop an algorithm $\mathcal A$ that reconstructs a dataset X_0' closely resembling the original forgetting set X_0 .

Evaluation Metric. The attack is considered successful if the attack algorithm \mathcal{A} reproduces the subsequent tokens exactly as they appear in the training set. The generated continuation is denoted by $\hat{x} = \mathcal{A}(\theta, \theta' \mid x_{\leq i})$, and the full set of extracted continuations over the dataset X is denoted as \hat{X} . By default, we treat the first half of each data sample as known and evaluate whether the attack algorithm can recover the remaining half.

We evaluate our method using the following two metrics:

- 1. **Rouge-L(R)**: Following previous work [23], we use Rouge-L [21] recall score (Rouge-L(R)) to measure the similarity between extracted continuations and the ground truth.
- 2. Average Extraction Success Rate (A-ESR $_{\tau}$): Inspired by prior work [3], we consider an extraction successful only if the generated sample is sufficiently similar to the ground truth. Formally, we define:

$$A-ESR_{\tau}(X_0, \widehat{X}) = \frac{1}{|X_0|} \sum_{i=1}^{|X_0|} Rouge-L(R)(X_0^{(i)}, \widehat{X}^{(i)}) \ge \tau.$$
 (1)

A threshold of $\tau=1.0$ indicates an exact match, while a $\tau<1.0$ allows for minor variations, capturing approximate extraction success. We measure A-ESR_{1.0} and A-ESR_{0.9} by default.

4 Proposed Method

4.1 Reversed Model Guidance

We illustrate the core idea of our method in Fig. 2. Building on prior work that successfully extracts finetuning data for diffusion models by guiding the transition from the model before fine-tuning to the model after fine-tuning [32], we view the unlearning process as the reverse of finetuning. We model this reversal as follows. Let the model before and after unlearning be denoted as θ and θ' , respectively. For exact unlearning, the only difference between these two models is whether the model has been trained on the forgetting set X_0 . We define $q(\cdot)$ as the ground truth probability of the forgetting set X_0 .

Given the unlearned model θ' , we assume a hypothetical process through which it relearns the distribution of

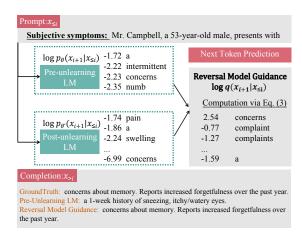


Figure 2: Visualization of reversed model guidance. We combine predictions from the pre- and post-unlearning models to approximate the forgotten distribution $q(x_{i+1}|x_{\leq i})$, resulting in a more effective extraction attack.

 X_0 , thereby approaching the original pre-unlearning model θ . This can be approximated by directly fine-tuning the model on the forgetting dataset X_0 . For any input $x_{\leq i}$, we then formulate the following parametric approximation for the next token prediction $p(x_{i+1}|x_{< i})$:

$$p_{\theta}(x_{i+1}|x_{\leq i}) \propto p_{\theta'}^{1-\lambda}(x_{i+1}|x_{\leq i})q^{\lambda}(x_{i+1}|x_{\leq i}),$$
 (2)

where λ is a coefficient, ranging from 0 to 1, that is related to the number of training iterations needed to adapt the model to the forgetting set X_0 . A higher λ corresponds to more training iterations, making the distribution $p_{\theta}(x)$ increasingly similar to the unlearned data distribution q(x).

Inspired by previous work applying classifier guidance in LLMs [29], we extend this concept to derive the log-probability form:

$$\log q(x_{i+1}|x_{\leq i}) = \log p_{\theta'}(x_{i+1}|x_{\leq i}) + w\left(\log p_{\theta}(x_{i+1}|x_{\leq i}) - \log p_{\theta'}(x_{i+1}|x_{\leq i})\right), \tag{3}$$

where $w = \frac{1}{\lambda}$ is the guidance scale, which is inversely proportional to the number of training iterations. With this model guidance, we simulate a "pseudo-predictor" $\log q(x_{i+1}|x_{\leq i})$ that steers the generation process toward high-probability regions within the unlearned data distribution q(x).

4.2 Token Filter Strategy

Directly using the log probability differences between two models can degrade generation quality and lead to incoherent or unnatural completions, as noted in previous work on contrastive decoding [19]. To mitigate this problem, we adopt the method in [19], which constrains token selection during decoding. For greedy decoding, this entails selecting the next token with the highest probability for the guided distribution $\log q$:

$$x_{\text{next}} = \arg\max_{v \in V'} \log q(v \mid x_{\leq i}), \tag{4}$$

but only within a constrained token set V' with high probability according to the pre-unlearning model θ :

$$V' = \{ v \in V \mid p_{\theta}(v \mid x_{\leq i}) \geq \gamma \max_{v \in V} p_{\theta}(v \mid x_{\leq i}) \}, \tag{5}$$

where V represents all possible tokens. The parameter γ controls the strictness of the candidate token filter. Intuitively, the pre-unlearning model retains residual knowledge of the unlearned dataset X_0 (otherwise, unlearning would be unnecessary). Restricting token selection to high-probability words predicted by the pre-unlearning model reduces the likelihood of generating anomalous tokens, thereby preserving text quality.

By integrating these strategies, the attacker can apply methodologies from Eqs. 4 and 5 to effectively generate text closely resembling the unlearning dataset X_0 .

5 Experiments

5.1 Experimental Setup

We evaluate unlearning methods on three datasets: the MUSE dataset [30], the TOFU dataset [23], and the WMDP dataset [18]. Following prior work [30, 23], we use Llama2-7B [31] and Phi-1.5 [20] as our base models. For each dataset, we first fine-tune the model on the full dataset to obtain the pre-unlearning checkpoint. We then apply exact unlearning by removing the forgetting set and re-fine-tuning the pretrained model on the remaining data.

Unless otherwise noted, we set the forgetting set size to 10% of the full dataset. For our method, the guidance scale w is set to 2.0 for Phi and 1.4 for Llama, and the constraint level γ is set to 10^{-5} by default. We analyze the impact of different fine-tuning iterations and forgetting set sizes in Sec. 5.3, and investigate the effect of varying hyper-parameters on the MUSE dataset in Sec. 5.4. Further details on training and dataset preparation are provided in Appendix Sec. A. We present additional results for our extraction method under LoRA fine-tuning in Appendix Sec. C, for larger LLMs in Appendix Sec. D, and for comparisons with other extraction attacks in Appendix Sec. E.

5.2 Main Comparison

To ensure fair comparison with previous work, we adopt a baseline attack that directly generates text from the given LLMs [25] before unlearning. Following prior studies [23, 30], we use greedy sampling by default, as it tends to exhibit higher memorization. We evaluate our method on multiple

Table 1: Comparison of our method and the baseline, which uses only the pre-unlearning model for extraction, across three datasets under various metrics. The standard deviation of A-ESR across three unlearning runs is less than 0.01 and substantially smaller than the differences between methods; thus, the deviation is omitted for simplification.

MUSE Dataset						
Phi-1.5 Llama2-7b						
	Rouge-L(R)↑	A -ESR $_{0.9}$ \uparrow	$A-ESR_{1.0}\uparrow$	Rouge-L(R)↑	A -ESR $_{0.9}$ \uparrow	$A-ESR_{1.0}\uparrow$
Post-unlearning Generation	0.296	0.006	0.004	0.212	0.014	0.013
Pre-unlearning Generation	0.473	0.114	0.101	0.675	0.424	0.384
Our Extraction	0.606	$0.249_{\substack{+118\%}}$	$0.224_{\uparrow 121\%}$	0.744	$0.496_{\uparrow 17.0\%}$	$0.438_{\uparrow 14.1\%}$
		TOF	U Dataset			
		Phi-1.5			Llama2-7b	
	Rouge-L(R)↑	A -ESR $_{0.9}$ \uparrow	A -ESR _{1.0} \uparrow	Rouge-L(R)↑	A -ESR $_{0.9}$ \uparrow	A -ESR $_{1.0}$ \uparrow
Post-unlearning Generation	0.437	0.007	0.005	0.420	0.012	0.010
Pre-unlearning Generation	0.566	0.100	0.070	0.588	0.185	0.093
Our Extraction	0.643	$0.202_{\uparrow 102\%}$	$0.120_{\uparrow 71.4\%}$	0.641	$0.218_{\uparrow 17.8\%}$	$0.133_{\uparrow 43.0\%}$
WMDP Dataset						
		Phi-1.5			Llama2-7b	
	Rouge-L(R) \uparrow	A -ESR $_{0.9}$ \uparrow	$A-ESR_{1.0}\uparrow$	Rouge-L(R) \uparrow	$A-ESR_{0.9}\uparrow$	$A-ESR_{1.0}\uparrow$
Post-unlearning Generation	0.278	0.011	0.009	0.222	0.006	0.006
Pre-unlearning Generation	0.429	0.079	0.069	0.313	0.062	0.050
Our Extraction	0.567	$0.218_{\uparrow 175\%}$	$0.192_{\uparrow 178\%}$	0.346	$0.087_{\uparrow 40.3\%}$	$0.075_{+50.0\%}$

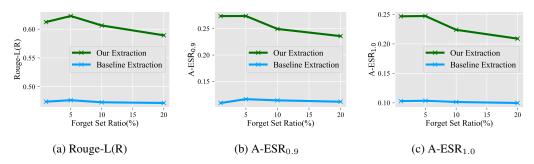


Figure 3: Comparison of our extraction method and the baseline on MUSE using Phi-1.5, evaluated at 3 epochs across different forgetting set ratios.

datasets (MUSE, TOFU, WMDP) using both Phi-1.5 and Llama2-7b, with 10% of the data designated as the forgetting set. As shown in Tab. 1, our method consistently achieves substantial improvements in extraction performance across all settings. Notably, the strict extraction accuracy (A-ESR($\tau=1.0$)) doubles in some cases and increases by at least $0.4\times$ in most settings, highlighting the effectiveness of our approach. Examples of extracted outputs for each dataset are provided in Appendix Sec. F.

5.3 Generalization

In this section, we further evaluate the applicability of our method across a broader range of scenarios, including varying forgetting set sizes and different numbers of training epochs. The former affects the overall difficulty of the unlearning task, as it determines how much the model's predictions are altered by the unlearning process, while the latter influences the extent to which the original model memorizes the forgetting set. We conduct experiments on the MUSE dataset using Phi-1.5, with the hyper-parameters fixed at w=2.0 and $\gamma=10^{-5}$.

Forgetting Set Size. As shown in Fig. 3, we observe that the forgetting set size has a relatively minor impact on extraction performance. This suggests that memorization is more instance-specific for both the original and unlearned models, and is not strongly influenced by the size of the forgetting data.

Training Epochs. As illustrated in Fig. 4, we find that with more training epochs—where the original model memorizes the forgetting set more extensively—the improvement from our method gradually diminishes. Our method is particularly effective when the model maintains a moderate level of memorization, which aligns with practical scenarios where models are trained for a moderate number of iterations to ensure good generalization while avoiding overfitting.

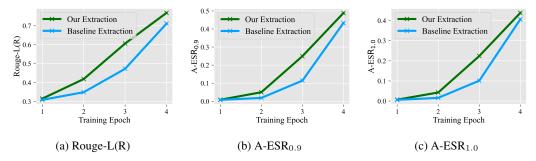


Figure 4: Comparison of our extraction method and the baseline on MUSE using Phi-1.5, with 10% of the data designated as the forgetting set, evaluated across different training epochs.

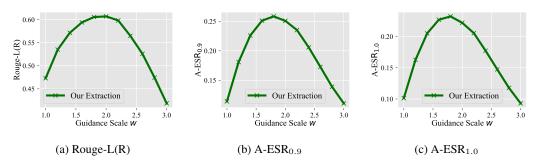


Figure 5: Extraction performance under different guidance scales w on MUSE using Phi-1.5, evaluated with a 10% forgetting set size.

5.4 Ablation Study

In this section, we experiment with the hyper-parameters in Eq. 3 and Eq. 5, including the guidance scale w and the token constraint strength γ . Experiments are conducted on the MUSE dataset with a 10% forgetting set size.

Guidance Scale w. The guidance scale w is the most critical hyper-parameter influencing extraction efficiency. Ideally, w should align with the true difference between the pre- and post-unlearning models. As shown in Fig. 5 and 6, w=2.0 works well for Phi-1.5, while w=1.4 is optimal for LLaMA2-7B.

We further investigate the optimal choice of w under different numbers of training epochs. As shown in Fig. 7, we observe that with larger training epochs—i.e., when the pre-unlearning model memorizes more—the optimal w becomes smaller. This observation aligns with the intuition derived from Eq. 2. According to Eq. 2, we assume an underlying fine-tuning process that transforms the post-unlearning model back into the pre-unlearning model. As the number of training epochs increases, a longer fine-tuning process would be needed, resulting in a larger λ , and consequently a smaller $w = \frac{1}{\lambda}$.

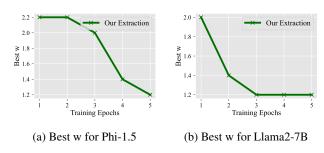


Figure 7: Optimal guidance scale w across different training epochs for Phi-1.5 and LLaMA2-7B. Experiments are conducted with a 10% forgetting set, and the best w is selected based on the highest Rouge-L(R) score. Results show that the optimal w decreases as training epochs increase.

Token Constraint Strength γ . In Eq. 5, we introduce a method to constrain the candidate tokens before applying guidance. We experiment with how γ influences extraction performance. As shown in Fig. 8, a moderate γ value between 10^{-3} and 10^{-5} generally improves performance. However, if γ is set too large, it interferes with the guidance signal and negatively impacts extraction effectiveness.

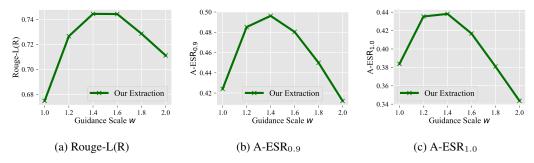


Figure 6: Extraction performance under different guidance scales w on MUSE using Llama2-7b, evaluated with a 10% forgetting set size.

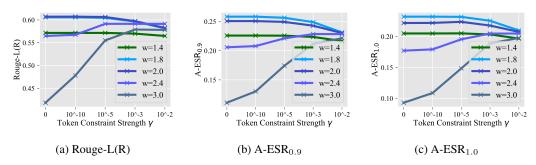


Figure 8: Extraction performance under different γ on MUSE using Phi-1.5, evaluated with a 10% forgetting set size.

5.5 Real-World Scenario Simulation: Extraction of Patient Information

We present a realistic and highly harmful scenario to illustrate the severity of our attack. Suppose a medical LLM has been fine-tuned on sensitive patient diagnostic records. We simulate this setting by constructing a dataset in real-world medical documentation formats [27], synthesized using Gemini 2.5 Pro. In this scenario, the attacker targets specific patients and may possess limited prior knowledge—such as the patient's name, birth date, or visit date.

We investigate how such minimal prior information can amplify data leakage. As shown in Table 2, our method yields a substantial improvement in extraction success rate, underscoring that such attacks can lead to severe privacy violations by effectively exposing a patient's sensitive information in real-world scenarios. Details on the medical dataset construc-

Table 2: Comparison of our method and the baseline on the medical dataset.

Medical Dataset						
Rouge-L(R) \uparrow A-ESR _{1.0} \uparrow						
Post-unlearning Generation	0.170	0				
Pre-unlearning Generation	0.320	0.140				
Our Extraction	0.457	0.210 _{↑50%}				

tion and an illustrative example are provided in Appendix Sec. B.

5.6 Extraction under Approximate Unlearning

We evaluate our extraction method under several approximate unlearning techniques [23, 35, 30]. Following our default setup, we fine-tune Phi-1.5 on the TOFU dataset for 3 epochs to obtain the pre-unlearning model, and experiment with a forgetting set that makes up 10% of the full dataset. For unlearning, we follow prior work [30], using a constant learning rate of 10^{-5} and stopping when the post-unlearning Rouge-L(R) score drops to or below that of exact unlearning. In our setting, this condition is consistently met after one epoch; further training leads to excessive utility degradation.

We evaluate the following representative approximate unlearning methods:

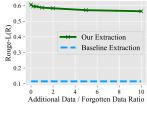
Table 3: Comparison of our extraction method with baselines under different unlearning methods on the TOFU dataset using Phi-1.5. Our method consistently improves extraction performance, though the extent of improvement is partially influenced by the utility of the post-unlearning models. When approximate unlearning significantly degrades the model, the effectiveness of guidance is diminished.

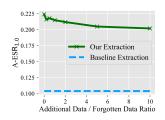
		Rouge-L(R)↑			A-ESR↑		Utility↑
	Post-Unlearning	Pre-Unlearning	Our Extraction	Post-Unlearning	Pre-Unlearning	Our Extraction	Post-Unlearning
Exact Unlearning (EU)	0.437	0.566	0.643	0.005	0.070	0.120 171.4%	0.567
GA	0.235	0.566	0.569	0.000	0.070	0.073	0.240
GA_{GD}	0.437	0.566	0.587	0.010	0.070	0.090	0.516
GAKL	0.243	0.566	0.571	0.002	0.070	0.075	0.253
NPO	0.272	0.566	0.579	0.003	0.070	0.080	0.282
NPO_{GD}	0.282	0.566	0.580	0.003	0.070	0.080	0.293
NPOKI	0.243	0.566	0.571	0.003	0.070	0.073	0.253

- Gradient Ascent (GA) [14, 13]: Applies gradient ascent on the cross-entropy loss to suppress the likelihood of the forget set. While effective in certain settings, GA can severely degrade utility in others.
- **Negative Preference Optimization (NPO)** [35]: Modifies the offline DPO objective to treat the forget set as negative preference data, encouraging low likelihood on it while remaining close to the original model.

To mitigate utility degradation, we incorporate two commonly used regularization strategies:

- Gradient Descent on the Retain Set (GD) [23]: Adds a standard cross-entropy loss on the retain set D_{retain} to maintain performance on non-forgotten data.
- KL Divergence Minimization (KL) [23]: Encourages the unlearned model's output distribution to remain close to that of the original model on inputs from the retain set.





(a) Rouge-L(R)

Figure 9: A-ESR_{1.0}

Figure 10: Effect of adding additional data as a defense on the MUSE dataset using Phi-1.5. The added data slightly reduces extraction performance.

We follow TOFU's default settings [23] for all approximate unlearning hyper-parameters. Following prior work [30, 23], utility is measured using the Rouge-L(R) score on the retain set.

For our extraction, we fix w=1.2 and $\gamma=10^{-5}$ across all approximate unlearning scenarios. As shown in Tab. 3, our method consistently improves extraction performance. However, the improvements are generally smaller than those observed under exact unlearning. We find that this reduction correlates with the utility of the post-unlearning model: as utility decreases, the benefit of guidance-based extraction also diminishes, as shown in Fig. 12. This suggests that approximate unlearning often sacrifices utility, which in turn distorts the guidance signal between the pre- and post-unlearning models, thereby reducing extraction effectiveness. This degradation is consistent with our observations in Sec. 5.7, where some defense strategies partially mitigate extraction risks but at the expense of model quality.

5.7 Possible Defense against the Attack

Adding Unrelated Data. Our extraction method relies on the difference between the pre- and post-unlearning models to capture the effect of removing the forgetting set. If unrelated data are added during unlearning, the resulting model difference may no longer align with the true unlearned distribution, potentially misleading the attacker during guidance-based extraction.

To evaluate whether this can serve as a viable defense, we conduct experiments on the MUSE dataset using a 10% forgetting set on Phi-1.5. We introduce auxiliary corpora from the WMDP dataset [18], which covers unrelated domains such as economics, law, physics, and cybersecurity, as additional data. During exact unlearning, we start from a pretrained LLM and fine-tune it on the full dataset excluding the forgetting set, augmented with varying amounts of the additional data.

As shown in Fig. 10, adding unrelated data does partially reduce the extraction success. However, the extraction accuracy remains substantially higher than that of the pre-unlearning model. Even with $10\times$ more unrelated data than the forgetting set—more than doubling the computational cost—the Rouge-L(R) and A-ESR metrics only exhibit a slight decline. This suggests that our extraction method is primarily instance-level and does not heavily rely on the model's overall conceptual knowledge, making it relatively insensitive to the introduction of additional unrelated data.

Noisy Gradient Updates. Inspired by Differential Privacy [5, 6], which provides theoretical guarantees against information leakage, we explore the use of DP-SGD [1] as a potential defense mechanism during exact unlearning. Specifically, we perturbed the updates with random noises before each gradient descent step. Intuitively, larger noise scales offer stronger privacy protection, but at the cost of reduced model utility.

We conduct experiments on the MUSE dataset with a 10% forgetting set using Phi-1.5, following the default setup. The only modification lies in the optimizer, where we inject Gaussian noise $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ with varying scales σ before each update. To quantify utility degradation, we follow prior work [30] and evaluate the Rouge-L(R) score on the retain set.

As shown in Fig. 11, increasing the noise scale consistently reduces the effectiveness of our extraction method. At sufficiently large noise levels (above 0.4), the extraction performance largely approaches that of the pre-unlearning model, indicating a partially effective defense. However, this comes at a significant cost: the model's utility on the retain set degrades severely, barely surpassing that of the original pretrained model. These findings suggest that while noisy gradient updates can serve as a partial defense, the trade-off between

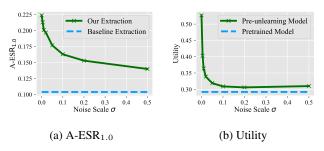


Figure 11: Effect of noisy gradient updates as a defense on the MUSE dataset using Phi-1.5. While large noise levels can partially mitigate our extraction attack, they also cause a substantial degradation in model utility.

privacy protection and model utility remains severe and undermines their practical viability.

6 Conclusion and Discussion

Most prior works on unlearning for LLMs focus solely on evaluating the privacy risk of the final unlearned model, without considering the implications of retaining access to earlier checkpoints or logits API. However, in many realistic scenarios—such as open-weight model releases or API deployments—there exists a practical risk that pre-unlearning models or logits may have been preemptively saved by an adversary. Our work shows that under such conditions, exact unlearning—widely regarded as the gold standard for data removal—can in a counterintuitive way introduce new privacy risks. By leveraging the differences between pre- and post-unlearning models through a guidance-based extraction method with token filtering, an adversary can significantly increase the leakage of the very content intended to be forgotten.

These findings reveal a previously overlooked but practical threat model. We urge the community to take this into account when designing and evaluating unlearning methods for LLMs. In particular, future techniques should offer privacy guarantees not only for the final model but also under adversarial access to its earlier states—only then can unlearning truly deliver on its intended privacy promises.

7 Acknowledgments

Zhiwei Steven Wu was in part supported by an NSF CAREER Award #2339775 and NSF Award #2232693.

References

- [1] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318, 2016.
- [2] M. Bertran, S. Tang, M. Kearns, J. H. Morgenstern, A. Roth, and S. Z. Wu. Reconstruction attacks on machine unlearning: Simple models are vulnerable. *Advances in Neural Information Processing Systems*, 37:104995–105016, 2024.
- [3] N. Carlini, F. Tramer, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. Brown, D. Song, U. Erlingsson, et al. Extracting training data from large language models. In 30th USENIX Security Symposium (USENIX Security 21), pages 2633–2650, 2021.
- [4] I. Consulting. General data protection regulation (gdpr), 2018. Accessed in April 2025.
- [5] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006*, New York, NY, USA, March 4-7, 2006. Proceedings 3, pages 265–284. Springer, 2006.
- [6] C. Dwork, A. Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends*® *in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- [7] R. Eldan and M. Russinovich. Who's harry potter? approximate unlearning for llms. 2023.
- [8] A. Ginart, M. Guan, G. Valiant, and J. Y. Zou. Making ai forget you: Data deletion in machine learning. *Advances in neural information processing systems*, 32, 2019.
- [9] Heidi Health. Heidi health ai medical scribe for global clinicians. https://www.heidihealth.com/, 2024.
- [10] E. J. Hu, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, W. Chen, et al. LoRA: Low-Rank Adaptation of Large Language Models. In *ICLR*, 2021.
- [11] H. Hu, S. Wang, T. Dong, and M. Xue. Learn what you want to unlearn: Unlearning inversion attacks against machine unlearning. In 2024 IEEE Symposium on Security and Privacy (SP), pages 3257–3275. IEEE, 2024.
- [12] S. Hu, Y. Fu, S. Wu, and V. Smith. Unlearning or obfuscating? jogging the memory of unlearned llms via benign relearning. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [13] G. Ilharco, M. T. Ribeiro, M. Wortsman, S. Gururangan, L. Schmidt, H. Hajishirzi, and A. Farhadi. Editing models with task arithmetic. *arXiv preprint arXiv:2212.04089*, 2022.
- [14] J. Jang, D. Yoon, S. Yang, S. Cha, M. Lee, L. Logeswaran, and M. Seo. Knowledge unlearning for mitigating privacy risks in language models. arXiv preprint arXiv:2210.01504, 2022.
- [15] J. Jia, J. Liu, P. Ram, Y. Yao, G. Liu, Y. Liu, P. Sharma, and S. Liu. Model sparsity can simplify machine unlearning. *Advances in Neural Information Processing Systems*, 36:51584–51605, 2023.
- [16] A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary, C. Bamford, D. S. Chaplot, D. d. l. Casas, E. B. Hanna, F. Bressand, et al. Mixtral of experts. arXiv preprint arXiv:2401.04088, 2024.
- [17] K. Kuo, A. Setlur, K. Srinivas, A. Raghunathan, and V. Smith. Exact unlearning of finetuning data via model merging at scale. *arXiv preprint arXiv:2504.04626*, 2025.
- [18] N. Li, A. Pan, A. Gopal, S. Yue, D. Berrios, A. Gatti, J. D. Li, A.-K. Dombrowski, S. Goel, L. Phan, et al. The wmdp benchmark: Measuring and reducing malicious use with unlearning. *arXiv preprint arXiv:2403.03218*, 2024.
- [19] X. L. Li, A. Holtzman, D. Fried, P. Liang, J. Eisner, T. Hashimoto, L. Zettlemoyer, and M. Lewis. Contrastive decoding: Open-ended text generation as optimization. In *The 61st Annual Meeting Of The Association For Computational Linguistics*, 2023.

- [20] Y. Li, S. Bubeck, R. Eldan, A. Del Giorno, S. Gunasekar, and Y. T. Lee. Textbooks are all you need ii: phi-1.5 technical report. arXiv preprint arXiv:2309.05463, 2023.
- [21] C.-Y. Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- [22] J. Łucki, B. Wei, Y. Huang, P. Henderson, F. Tramèr, and J. Rando. An adversarial perspective on machine unlearning for ai safety. *arXiv preprint arXiv:2409.18025*, 2024.
- [23] P. Maini, Z. Feng, A. Schwarzschild, Z. C. Lipton, and J. Z. Kolter. Tofu: A task of fictitious unlearning for llms. arXiv preprint arXiv:2401.06121, 2024.
- [24] K. K. Nakka, A. Frikha, R. Mendes, X. Jiang, and X. Zhou. Pii-compass: Guiding Ilm training data extraction prompts towards the target pii via grounding. arXiv preprint arXiv:2407.02943, 2024.
- [25] M. Nasr, N. Carlini, J. Hayase, M. Jagielski, A. F. Cooper, D. Ippolito, C. A. Choquette-Choo, E. Wallace, F. Tramèr, and K. Lee. Scalable extraction of training data from (production) language models. *arXiv preprint arXiv:2311.17035*, 2023.
- [26] S. of California Department of Justice. California consumer privacy act (ccpa), 2018. Accessed in April 2025.
- [27] V. Podder, V. Lew, and S. Ghassemzadeh. SOAP Notes. https://www.ncbi.nlm.nih.gov/books/NBK482263/, 2023. StatPearls [Internet]. Treasure Island (FL): StatPearls Publishing; Updated 2023 Aug 28.
- [28] H. Qiu, Y. Wang, Y. Xu, L. Cui, and Z. Shen. Fedcio: Efficient exact federated unlearning with clustering, isolation, and one-shot aggregation. In 2023 IEEE International Conference on Big Data (BigData), pages 5559–5568. IEEE, 2023.
- [29] G. Sanchez, A. Spangher, H. Fan, E. Levi, and S. Biderman. Stay on topic with classifier-free guidance. In *Forty-first International Conference on Machine Learning*, 2024.
- [30] W. Shi, J. Lee, Y. Huang, S. Malladi, J. Zhao, A. Holtzman, D. Liu, L. Zettlemoyer, N. A. Smith, and C. Zhang. Muse: Machine unlearning six-way evaluation for language models. *arXiv* preprint arXiv:2407.06460, 2024.
- [31] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv* preprint arXiv:2302.13971, 2023.
- [32] X. Wu, J. Zhang, and S. Wu. Revealing the unseen: Guiding personalized diffusion models to expose training data. *arXiv* preprint arXiv:2410.03039, 2024.
- [33] X. Xia, Z. Wang, R. Sun, B. Liu, I. Khalil, and M. Xue. Edge unlearning is not" on edge"! an adaptive exact unlearning system on resource-constrained devices. *arXiv* preprint *arXiv*:2410.10128, 2024.
- [34] Z. Xiong, W. Li, Y. Li, and Z. Cai. Exact-fun: an exact and efficient federated unlearning approach. In 2023 IEEE International Conference on Data Mining (ICDM), pages 1439–1444. IEEE, 2023.
- [35] R. Zhang, L. Lin, Y. Bai, and S. Mei. Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv e-prints*, pages arXiv–2404, 2024.
- [36] Z. Zhang, J. Wen, and M. Huang. Ethicist: Targeted training data extraction through loss smoothed soft prompting and calibrated confidence estimation. *arXiv* preprint *arXiv*:2307.04401, 2023.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction highlight a new threat: that exact unlearning—previously regarded as the gold standard for data removal—can introduce new privacy risks. This claim is the central focus of the paper and is supported by the proposed method in Sec. 4 and a set of experiments that design and evaluate a data extraction attack leveraging both pre- and post-unlearning models in Sec. 5.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
 contributions made in the paper and important assumptions and limitations. A No or
 NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The paper discusses the limitations in Appendix Sec. G.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper includes heuristic modeling to provide intuitive guidance for the proposed method in Sec. 4, but it does not involve formal theoretical assumptions or proofs.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The main experimental procedures are outlined in Sec. 5, with additional implementation details provided in Appendix Sec. A, ensuring reproducibility of the results relevant to the core claims.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in

some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The code is provided via an link shown at abstract. The datasets used are primarily open-source benchmarks. The main experimental results are presented in Sec. 5, with reproduction details provided in Appendix Sec. A.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The main experimental settings are described in Sec. 5, with additional details provided in Appendix Sec. A. Hyper-parameters and their selection process are explained in Sec. 5.4, ensuring clarity on how they were chosen.

Guidelines:

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Statistical significance is addressed in Tab. 1. The variation of A-ESR across three unlearning runs is minimal (standard deviation < 0.01), so error bars are omitted for clarity. The improvements of our extraction method over the baseline are substantially larger than this variance, indicating the results are statistically robust.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Information available in Sec. A

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We have reviewed the NeurIPS Code of Ethics and ensured that the research complies with all relevant ethical guidelines.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The paper discusses the broader impacts in Appendix Sec. G.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper does not center on releasing any new datasets or models, and thus no additional safeguards are required.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Ouestion: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We use multiple benchmark datasets, all of which are properly cited in Sec. 5, with detailed descriptions provided in Appendix Sec. A.

Guidelines:

• The answer NA means that the paper does not use existing assets.

- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: LLMs are used during the construction of the medical dataset, as described in Sec. 5.5, but they are not part of the core methodology or contribution of the paper.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

A Experiment Details

All experiments are conducted using two NVIDIA A100 GPUs.

A.1 Training Details

Following prior works [30, 23], we begin with the pre-trained LLMs, LLaMA2-7B and Phi-1.5. To obtain the target (pre-unlearning) model with moderate memorization of the training data, we fine-tune LLaMA2-7B for 2 epochs and Phi-1.5 for 3 epochs on the full dataset, using a constant learning rate of 10^{-5} . This setup reflects a realistic scenario: a well-tuned model should neither memorize excessively—compromising generalization—nor memorize too little, which would eliminate the need for unlearning in the first place.

To simulate exact unlearning, we train a second model from scratch with the same configurations, again starting from the pre-trained weights but excluding the designated forgetting set. This yields the post-unlearning model for our experiments.

A.2 Dataset Preparation

We experiment on the following benchmark datasets: MUSE, TOFU, and WMDP.

- MUSE [30]: We use the MUSE-News dataset, which consists of BBC news articles collected after August 2023. The dataset is split into two disjoint subsets: $\mathcal{D}_{\text{forget}}$ and $\mathcal{D}_{\text{retain}}$, containing 0.8M and 1.6M tokens, respectively. For a k% forgetting set, we randomly select passages from $\mathcal{D}_{\text{forget}}$ until the total number of selected tokens reaches $2.4\text{M} \times k\%$. The prefix known to the attacker is the first half of each sentence.
- **TOFU** [23]: We use the full TOFU dataset, which consists entirely of fictitious author biographies synthesized by GPT-4. To construct the forgetting set, we randomly sample question-answer pairs and treat the remaining data as the retaining set. The prefix known to the attacker is the question part.
- WMDP [18]: We use a subset of bio-retain-corpus from WMDP, comprising a collection of PubMed papers that span various categories within general biology. This subset contains a total of 5.3k sentences. We randomly sample sentences from this subset to form the forgetting set, with the remainder serving as the retaining set. We simulate the attacker's prior knowledge by providing access to the first half of each sentence as a prefix.

B Medical Dataset Experiment Details

To simulate real-world medical data, we design our medical dataset using the Subjective, Objective, Assessment and Plan (SOAP) note [27] as a template. SOAP notes are a widely adopted method for healthcare providers to document patient encounters in a structured and organized manner. To ensure a comprehensive structure, we utilized the SOAP note template from Heidi Health [9], a medical AI company which offers SOAP note templates provided by specialists from the medical industry.

We format our dataset into JSON with the following keys, "client name", "date of birth", "date", "subjective", "objective", "assessment" and "plan". For the generation process, we employ Gemini 2.5 Pro using a specialized prompt:

I would like to generate synthetic medical data for machine learning purposes. Specifically, I would use SOAP notes as the data type. Below is a note template you need to follow, which has client name, date of birth, date, as well as subjective, objective, assessment, and plan. The template is just for you to refer, you do not need to generate each line of the template. Instead, only several lines for each of the SOAP is enough, try not to be too tedious for each record. For each record, please generate with a PII (client name, date of birthday), one person per record. client name: [name]

date of birth: [birthday date]

date: [visiting date]

Subjective:

[Description of symptoms, onset of symptoms, location of symptoms, duration of symptoms, characteristics of symptoms, alleviating or aggravating factors, timing, and severity]

[Current medications and response to treatment] (write this section in narrative form. Write in full sentences and do not include any bullet points)

[Any side effects experienced] (write this section in narrative form. Write in full sentences and do not include any bullet points)

[Non-pharmacological interventions tried] (write this section in narrative form. Write in full sentences and do not include any bullet points)

[Description of any related lifestyle factors] (write this section in narrative form . Write in full sentences and do not include any bullet points)

[Patient's experience and management of symptoms] (write this section in narrative form. Write in full sentences and do not include any bullet points)

[Any recent changes in symptoms or condition] (write this section in narrative form. Write in full sentences and do not include any bullet points)

[Any pertinent positive or pertinent negatives in review of systems] (write this section in narrative form. Write in full sentences and do not include any bullet points)

Objective:

Vital Signs:Blood Pressure: [blood pressure reading] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)

Heart Rate: [heart rate reading] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)
Respiratory Rate: [respiratory rate reading] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)
Temperature: [temperature reading] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)
Oxygen Saturation: [oxygen saturation reading] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)
General Appearance: [general appearance description] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)

HEENT: [head, eyes, ears, nose, throat findings] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)

Neck: [neck findings] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)

Cardiovascular: [cardiovascular findings] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)
Respiratory: [respiratory findings] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)
Abdomen: [abdominal findings] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)
Musculoskeletal: [musculoskeletal findings] (only include if explicitly mentioned in

the transcript, contextual notes or clinical note, otherwise leave blank.) Neurological: [neurological findings] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)

Skin: [skin findings] (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank.)

Assessment:

[Likely diagnosis]

[Differential diagnosis (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank)]

Diagnostic Tests: (only include if explicitly mentioned other skip section) [Investigations and tests planned (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank)]
Plan:

[Treatment planned for Issue 1 (only include if explicitly mentioned in the transcript, contextual notes or clinical note, otherwise leave blank)]
[Relevant referrals for Issue 1 (only include if explicitly mentioned- [Likely diagnosis for Issue 1 (condition name only)]

(Never come up with your own patient details, assessment, diagnosis, interventions, evaluation or plan for continuing care - use only the transcript, contextual notes, or clinical note as a reference for the information included in your note. If any

information related to a placeholder has not been explicitly mentioned in the transcript, contextual notes, or clinical note, you must not state the information has not been explicitly mentioned in your output, just leave the relevant placeholder or section blank).

Then, here is an example of SOAP fields for you to refer: Subjective:

The patient, a 52-year-old male, presents with a new rash on his back and arms, which he has noticed for the past two weeks. He describes the rash as "itchy and red ," and mentions that it seems to be getting worse despite over-the-counter anti-itch creams. The patient denies any fever, joint pain, or recent exposure to new soaps or detergents.

Objective:

Appearance: The patient appears well-nourished and in no acute distress. Skin: Exam reveals erythematous, scaly plaques on the back and arms. There is evidence of excoriation due to itching. No signs of systemic involvement. Lesions: Lesions are well-defined, with some areas showing mild papules. No signs of pustules or ulcers.

Other Systems: Vital signs are within normal limits. No lymphadenopathy noted. Assessment:

The presentation is consistent with psoriasis, characterized by itchy, scaly plaques . The absence of systemic symptoms and well-defined lesions supports this diagnosis. Differential diagnoses include eczema or fungal infection, but these are less likely given the clinical presentation. Plan:

Initiate topical treatment with high-potency corticosteroids to reduce inflammation and itching.

Recommend emollients to improve skin hydration and prevent dryness.

Educate the patient on the nature of psoriasis, including triggers and management strategies.

Suggest lifestyle modifications such as stress management and dietary adjustments to potentially improve symptoms.

You should add PII in front of the SOAP. Please generate 10 records for me, in json format, with "client name, date of birth, date, as well as subjective, objective, assessment, and plan" as keys.

To maintain data quality and prevent degradation during large-scale generation, we created the dataset in batches of 50 records each, producing 1,000 records in total. We replaced duplicate client names with unique ones, as there would be a low probability of identical names appearing in a real-world sample of this size. The generated data covered a diverse range of medical conditions to ensure a representative sample of real-world clinical scenarios. We randomly sample 100 records as the forgetting set, with the remaining 900 records serving as the retaining set.

For illustration purposes, an example of the generated records is provided below:

"client name": "Noah Garcia",

```
"date of birth": "2012-07-22",
"date": "2025-05-18",
"subjective": "Parent reports Noah, a 12-year-old male, has had intermittent
abdominal pain for the past month. Pain is periumbilical, crampy, occurs 1-2 times
per week, lasting 30-60 minutes. No clear relation to food. No fever, vomiting,
diarrhea, or weight loss. Appetite is normal. School attendance is unaffected. He
takes no medications. Parent has tried giving children's Tylenol during episodes
with little effect. Parent is worried about the recurrence.",
"objective": "Vital Signs: Normal for age. Abdomen: Soft, non-tender, non-distended.
Bowel sounds normal. No masses palpated. Growth chart parameters are normal.",
```

characteristics, and lack of red flag symptoms.", "plan": "Reassure parent and child about functional nature. Discuss potential triggers (stress, diet). Recommend keeping a pain and stool diary. Encourage highfiber diet and adequate fluids. Advise follow-up if pain changes pattern, becomes

"assessment": "Recurrent abdominal pain, likely functional abdominal pain given age,

severe, or if red flag symptoms (weight loss, vomiting, blood in stool) develop."

To capture the complex structure of these long medical records, we fine-tune the model for 11 epochs, while keeping all other settings unchanged in the experiment.

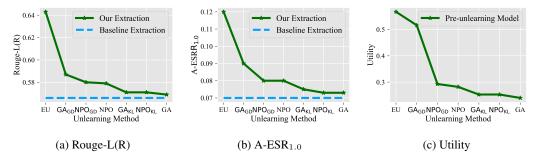


Figure 12: Comparison of our extraction method against the baseline under various unlearning methods. EU refers to exact unlearning. In some cases, the effectiveness of our method weakens—primarily due to reduced model utility, which distorts the guidance between the pre- and post-unlearning models.

Table 4: Comparison of our extraction method with baselines under LoRA fine-tuning on Phi-1.5 on TOFU dataset. The improvements are consistent across metrics.

Extraction Methods	Rouge- $L(R)\uparrow$	A-ESR _{0.9} ↑	A-ESR _{1.0} \uparrow
Post-Unlearning Generation	0.412	0.010	0.003
Pre-Unlearning Generation	0.544	0.070	0.055
Ours	0.614	0.160	0.108

C Extraction Results under LoRA Fine-tuning

In practice, parameter-efficient fine-tuning methods such as LoRA [10] are increasingly popular. We therefore examine whether our attack remains effective when the pre-unlearning model is LoRA-fine-tuned on TOFU. Specifically, we provide results on Phi-1.5 fine-tuned with LoRA on TOFU for five epochs using a constant learning rate of 2e-4 and a 10% forgetting set, while keeping all other settings consistent with Sec. 5.2. Then, we employ a guidance scale w=1.6 for our extraction, increasing the extraction rate. The improvements are consistent, as shown in Tab. 4.

D Extraction Results on Larger LLMs

We further experiment with a larger model, Mixtral-8x7B [16], which has substantially more parameters, employs a MoE architecture, and is instruction-tuned, making it sufficiently distinct from our default models (LLaMA and Phi) to better evaluate generalization. Specifically, we report results on Mixtral-8x7B-Instruct-v0.1², fine-tuned with LoRA on TOFU for two epochs using a cosine learning rate schedule with a base learning rate of 1e-4 and a 10% forgetting set (all other settings follow Sec. 5.2). For extraction, we apply a guidance scale of w=2. The improvements remain consistent, as shown in Tab. 5.

E Comparison with Other Extraction Attacks

To further contextualize our method within the broader landscape of extraction attacks, we compare it against ETHICIST [36], a representative approach for extracting memorized training data. ETHICIST operates by tuning soft prompt embeddings and applying loss smoothing with calibrated confidence estimation. We conduct experiments on the TOFU dataset using Phi-1.5 under the 10% forgetting setting.

Notably, ETHICIST assumes a stronger threat model that the attacker has access to part of the model's training data. To align with this setting, we assume the attacker has access to half of the retained set in TOFU (while remaining blind to the target extraction set, i.e., the forgetting set), and use it to

²https://huggingface.co/mistralai/Mixtral-8x7B-Instruct-v0.1

Table 5: Comparison of our extraction method with baselines under LoRA fine-tuning on Mixtral-8x7B-Instruct-v0.1 on TOFU dataset. The improvements are consistent across metrics.

Extraction Methods	Rouge- $L(R)\uparrow$	A-ESR _{0.9} \uparrow	A-ESR _{1.0} \uparrow
Post-Unlearning Generation	0.372	0.015	0.013
Pre-Unlearning Generation	0.465	0.040	0.018
Ours	0.537	0.103	0.055

Table 6: Comparison of our extraction method with ETHICIST under default fine-tuning on the TOFU dataset. Our method outperforms ETHICIST by a clear margin, and combining the two further improves performance.

Extraction Methods	Rouge-L(R)↑	A-ESR _{0.9} ↑	A-ESR _{1.0} ↑
Pre-Unlearning Generation	0.566	0.100	0.070
ETHICIST	0.570	0.118	0.073
Ours	0.643	0.202	0.120
Ours + ETHICIST	0.652	0.213	0.153

train the soft prompt. The hyper-parameters of ETHICIST follow the configuration provided in the original paper's open-sourced code.

As shown in Tab. 6, experimental results show that ETHICIST alone can partially improve the extraction rate. However, since ETHICIST is designed to extract training data from a single model, it is complementary rather than conflicting with our method, which leverages the difference between pre- and post-unlearning models. In fact, combining ETHICIST with our guidance-based approach leads to even stronger extraction performance.

F Visualization

In Figs. 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, and 24, we present several examples under our default setting using Phi-1.5, with our method applied using the default hyper-parameters (w=2.0, $\gamma=10^{-5}$). For each dataset, we include examples where both our method and the baseline fail, where both succeed, and intermediate cases where the baseline fails but our method successfully improves extraction.

G Limitations and Broader Impact

In this paper, we show that exact unlearning—originally intended to improve model safety—can, in fact, introduce new privacy risks. Our method relies on access to weights or logits api from both the pre- and post-unlearning models. While we justify this assumption using a realistic medical dataset, there are cases where pre-unlearning checkpoints or logits may not be available, such as in closed-source settings or when attackers fail to pre-save sufficient outputs. This limits the general applicability of our method. Future work may explore leveraging public model outputs or general-purpose knowledge priors, as suggested in prior work [2].

Our extraction method reveals a privacy risk of exact unlearning that, in principle, could be exploited in practice. However, as with other papers that focus on attacks, our goal is not to promote misuse, but to highlight a potential vulnerability before it leads to real-world consequences. By identifying this risk early, we hope to encourage more cautious use of exact unlearning and to motivate the community to proactively develop stronger defense mechanisms.

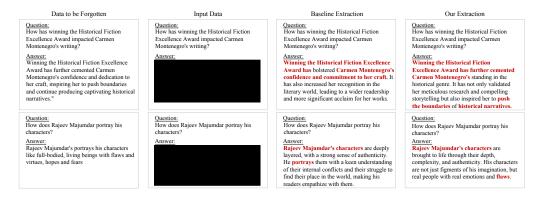


Figure 13: Examples from the TOFU dataset illustrating hard extraction cases where both our method and the baseline fail.

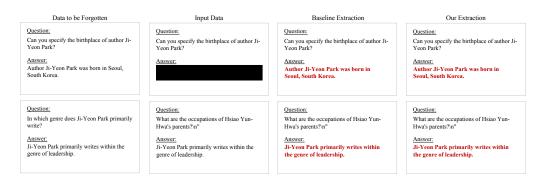


Figure 14: Examples from the TOFU dataset illustrating easy extraction cases where both our method and the baseline mostly succeed.

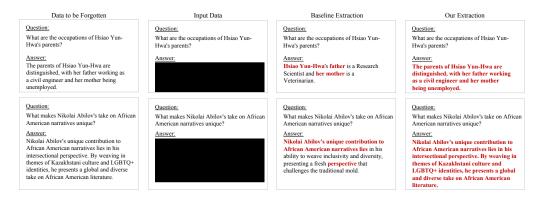


Figure 15: Examples from the TOFU dataset showing cases of intermediate extraction difficulty, where the baseline fails, but our method successfully recovers most of the target information. These cases highlight the improvement brought by the proposed method.

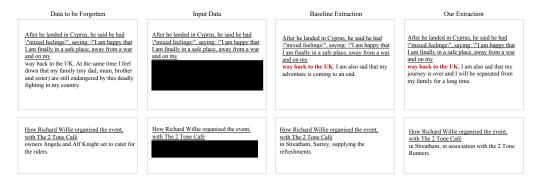


Figure 16: Examples from the MUSE dataset illustrating hard extraction cases where both our method and the baseline fail.

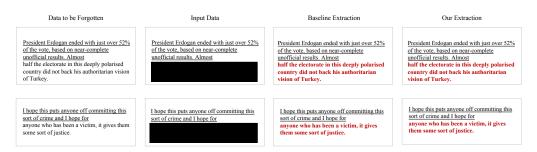


Figure 17: Examples from the MUSE dataset illustrating easy extraction cases where both our method and the baseline mostly succeed.

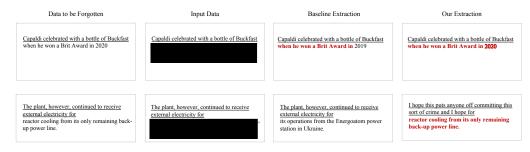


Figure 18: Examples from the MUSE dataset showing cases of intermediate extraction difficulty, where the baseline fails, but our method successfully recovers most of the target information. These cases highlight the improvement brought by the proposed method.

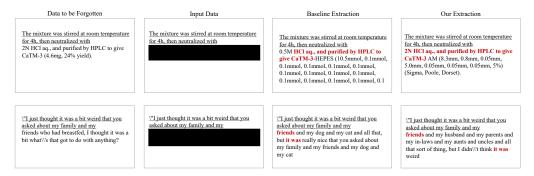


Figure 19: Examples from the WMDP dataset illustrating hard extraction cases where both our method and the baseline fail.

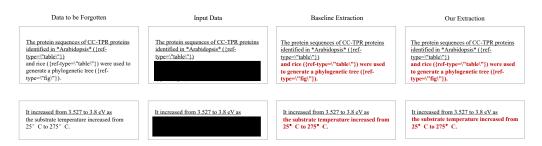


Figure 20: Examples from the WMDP dataset illustrating easy extraction cases where both our method and the baseline mostly succeed.

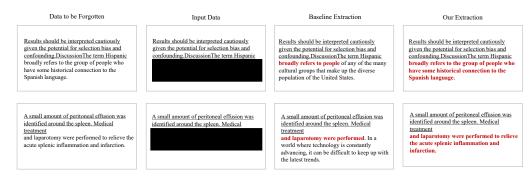


Figure 21: Examples from the WMDP dataset showing cases of intermediate extraction difficulty, where the baseline fails, but our method successfully recovers most of the target information. These cases highlight the improvement brought by the proposed method.



Figure 22: Examples from the medical dataset illustrating hard extraction cases where both our method and the baseline fail.

Patient data to be forgotten unlearning checkpoints client name: Hannah Gates date of birth:
1993-01-23 date: 2025-08-29
subjective: Ms. Gates, 32, for follow-up of
an incidental finding of a 3cm simple
ovarian cyst on the left ovary found on
pelvic ultrasound done for unrelated reasons
3 months ago. She is asymptomatic. No
abdominal pain, bloating, or changes in
menstrual cycle. She is using oral
contraceptives. No family history of ovarian
caneer, objective: Vital Siens: Normal. client name: Hannah Gates date of birth:
1993-01-23 date: 2025-08-29
subjective: Ms. Gates, 32, for follow-up of
an incidental finding of a Sem simple
ovarian cyst on the left ovary found on
pelvic ultrasound done for unrelated
reasons 3 months ago. She is
asymptomatic. No abdominal pain,
bloatine, or changes in menstrual cycle. client name: Hannah Gates date of birth: 1993-01-23 date: 2025-08-29 subjective: Ms. Gates, 32, for follow-up of an incidental finding of a 3cm simple ovarian cyst on the left ovary found on client name: Hannah Gates date of birth: pelvic ultrasound done for unrelated reasons 3 months ago. She is asymptomatic. No abdominal pain, bloating, or changes in menstrual cycle. She is using oral contraceptives. No family bloating, or changes in menstrual cycle. She is using oral contraceptives. No family cancer. objective: Vital Signs: Normal. history of ovarian cancer. objective: Vital history of ovarian cancer. objective: Vital Abdomen: Soft, non-tender. Pelvic exam (if performed, or defer if not indicated for Signs: Normal. Abdomen: Soft, non-Signs: Normal. Abdomen: Soft, nonperformed, or defer if not indicated for asymptomatic follow-up): Unremarkable Review of prior ultrasound report confirming simple cyst characteristics. assessment: Likely diagnosis: Asymptomatic simple ovarian cyst, left, likely functional. Low risk of malignancy in premenopausal woman with simple cyst <5cm. plan: Reassure patient that simple cysts are common and usually benign, often resolving spontaneously. Given size and simple nature, recommend follow-up pelvic ultrasound in 3-6 months to ensure stability or resolution. Continue current contraception. Advise to return if new symptoms develop (e.g., pelvic pain, bloating). No other intervention needed at this time. tender. Pelvic exam (if performed, or defer tender. Pelvic exam (if performed, or defer if not indicated for asymptomatic follow-up): Unremarkable. Review of prior ultrasound report confirming simple cyst characteristics. assessment: Likely if not indicated for asymptomatic follow-up): Unremarkable. Review of prior ultrasound report confirming simp characteristics, assessment: Likely characteristics, assessment: Likely diagnosis: Asymptomatic simple ovarian cyst. Ieft, likely functional. Low risk of malignancy in premenopausal woman with simple cyst <5cm, plan: Reassure patient that simple cysts are common and usually benign, often resolving spontaneously. Given size and simple nature, recomment follow-up pelvic ultrasound in 3-6 months to ensure stability or resolution. Continue current contraception. Advise to return if new symptoms develop (e.g., pelvic pain, characteristics, assessment: Likely diagnosis: Asymptomatic simple ovarian eyst, left, likely functional. Low risk of malignancy in premenopausal woman with simple eyst <5cm, plan: Reassure patient that simple eysts are common and usually benigh, often resolving spontaneously. Given size and simple nature, recommend follow-up pelvic ultrasound in 3-6 months to ensure stability or resolution. Continue current contraception. Advise to return if new symptoms develop (e.g., pelvic pain, at this time. symptoms develop (e.g., pelvic pain, bloating). No other intervention needed at symptoms develop (e.g., pelvic pain, bloating). No other intervention needed at client name: Aurelia Lambert date of birth: 1997-11-04 date: 2026-06-05 subjective: Ms. Lambert is a 28-year-old female (MSM context - uses she/her pronouns) presenting for discussion about HIV prevention. Reports multiple male partners, inconsistent condom use. Aware of PrEP, wants to start, objective: Sexual history confirms substantial risk. HIV Ag/Ab test negative today. HBV immune (prior vaccine). Renal function normal. assessment: Appropriate candidate for HIV PrEP based on risk factors, plan: Initiate PrEP with daily oral F/TAF (Descoy). Counsed on importance of adherence for effectiveness. Discuss potential side effects. Schedule 3-month follow-up for repeat HIV testing, STI screening (GCC/TIS)philis), renal function monitoring, adherence check, and prescription refill. Reinforce safer sex practices including consistent condom use for STI nevention client name: Aurelia Lambert date of birth: 1997-11-04 date: 2026-06-05 subjective: Ms. Lambert is a 28-year-old female (MSM context - uses sheher pronouns) presenting for discussion about HIV prevention. Reports multiple male partners, inconsistent condom use. Aware of PrEP, wants to start, objective: Sexual history confirms substantial risk. HIV Ag/Ab test negative today. HBV immune (prior vaccine). Renal function normal. assessment: Appropriate candidate for HIV PrEP based on risk factors, plan: Initiate PrEP with daily oral F/TAF (Descoy). Counsed on importance of adherence for effectiveness. Discuss potential side effects. Schedule 3-month follow-up for repeat HIV testing, STI screening (GCC/TS/Sphilis), renal function monitoring, adherence check, and prescription refill. Reinforce safer sex. client name: Aurelia Lambert date of birth: 1997-11-04 date: 2026-06-05 subjective: Ms. Lambert is a 28-year-old female (MSM context - uses she/her pronouns) presenting for discussion about client name: Aurelia Lambert date of birth: 1997-11-04 date: 2026-06-05 pronous) presenting for discussion about HIV prevention. Reports multiple male partners, inconsistent condom use. Aware of PrEP, wants to start, objective: Sexual history confirms substantial risk. HIV Ag/Ab test negative today. HBV immune (prior vaccine). Renal function normal. assessment: Appropriate candidate for HIV PrEP based on risk factors. plan: initiate PrEP with daily oral FTAF (Descovy). Counsel on importance of adherence for effectiveness. Discuss potential side effects. Schedule 3. month follow-up for repeat HIV testing, STI screening (GC/CT/Syphilis), renal function monitoring, adherence check, and monitoring, adherence check, and prescription refill. Reinforce safer sex function monitoring, adherence check, and prescription refill. Reinforce safer sex practices including consistent condom use for STI prevention. practices including consistent condom use for STI prevention. practices including consistent condom use for STI prevention.

Baseline Extraction

using pre-unlearning checkpoint

Our Extraction: combining pre- and post-

Adversary's side information

about the patient

Figure 23: Examples from the medical dataset illustrating easy extraction cases where both our method and the baseline mostly succeed.

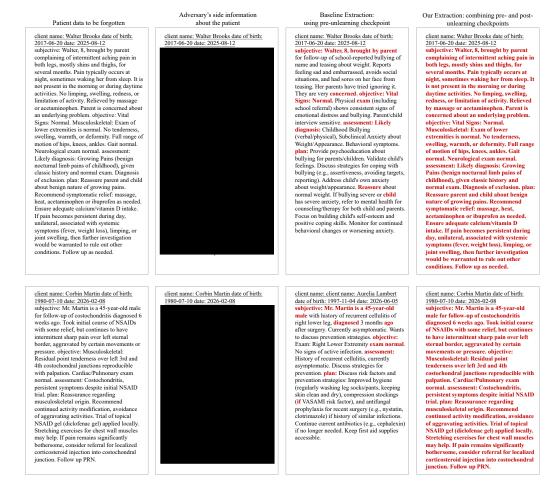


Figure 24: Examples from the medical dataset showing cases of intermediate extraction difficulty, where the baseline fails, but our method successfully recovers most of the target information. These cases highlight the improvement brought by the proposed method.