

UNRE: ZERO-SHOT LLM UNLEARNING VIA DYNAMIC CONTEXTUAL RETRIEVAL

Anonymous authors

Paper under double-blind review

ABSTRACT

Inference-time machine unlearning with only the forget data, also known as zero-shot unlearning, is becoming increasingly important for bias mitigation, privacy preservation, copyright protection, etc. Most approaches in this domain focused on query updating, decoder modification, offline module training, or reverse-generation by the forget data. Recent works found that providing offline-prepared contexts can realize in-context unlearning. However, leveraging dynamic context (conditioned on real-time queries) to achieve zero-shot unlearning has not yet been explored, which has the potential to enforce context unlearning while preserving the performance of the original LLM. In this paper, we propose UNRE, a novel unlearning framework for LLMs that employs dynamic contextual retrieval from retrieval-augmented generation (RAG) while only leveraging the forget data. Specifically, UNRE dynamically updates contexts to guide the unlearning process in a zero-shot unlearning setting. During the inference, the user query is first leveraged for online membership inference to identify a query-specific forget set. Using this set, UNRE refines the embeddings of the retrieved chunks via gradient descent, producing adaptive contexts that steer the LLM toward a query-specific unlearned distribution. We evaluate UNRE on multiple unlearning benchmarks and show that UNRE not only outperforms existing zero-shot and context-based unlearning approaches, but also better preserves the original model performance.

1 INTRODUCTION

Machine unlearning is the process of revoking or forgetting data embedded in the memory of a pre-trained model (Bourtoule et al., 2021). Unlike catastrophic forgetting (Goodfellow et al., 2013), which arises unintentionally during training, machine unlearning aims to deliberately and controllably erase specific knowledge from a model. Effective unlearning is critical for building trustworthy large language models (LLMs), as it enables the removal of harmful responses (Yao et al., 2024a; Li et al., 2024; Barrett et al., 2023), copyrighted content (Dou et al., 2025; Chen et al., 2023), societal biases (Motoki et al., 2024; Yu et al., 2023), hallucinations (Yao et al., 2024a), and supports timely safety alignment (Song et al., 2025). Traditional machine unlearning methods can be categorized into targeted and untargeted approaches (Yuan et al., 2025). These methods typically require not only a *forget set*—the data to be removed from the model—but also either a reference model (Ji et al., 2024) or a *retain set*, i.e., the original training data excluding the *forget set*. The retain set can be constructed through membership inference (Shokri et al., 2017), reverse generation from the *forget set* (Pawelczyk et al., 2024), and related techniques. However, since the retain set is often unavailable in real-world scenarios (Li et al., 2024), recent works such as FLAT (Wang et al., 2025b) have been proposed to enable unlearning using only the *forget data*. Zero-shot unlearning has emerged as a scenario where the source training data is unavailable (Chundawat et al., 2023; Foster et al., 2024; Chen et al., 2025; Ahmed et al., 2025); instead, the method only requires the forget request data.

LLM unlearning targets the removal of knowledge in a designated *forget set* while preserving the model performance on other tasks (Wang et al., 2025b). Beyond data-based approaches described above, other methods include model-based unlearning, which relies on fine-tuning (Yao et al., 2024a) or training specific modules (Bhaila et al., 2025), and input-based unlearning (Liu et al., 2024a; Pawelczyk et al., 2024). Input-based methods (Liu et al., 2024a) achieve unlearning by modifying the prompt (e.g., gradient-based updates of prompt embeddings (Bhaila et al., 2025; Liu et al.,

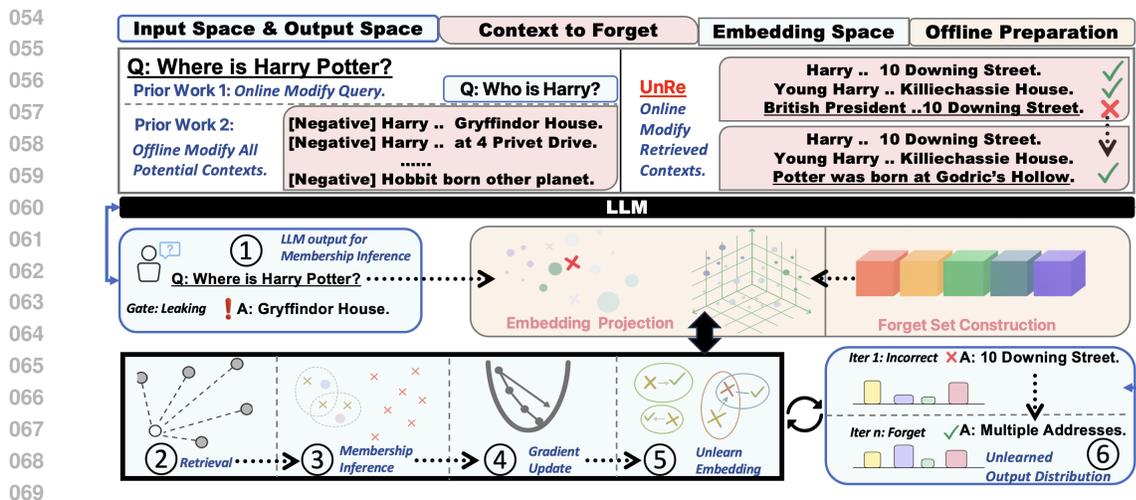


Figure 1: **Upper section: visualizing dynamic modification of the context by UNRE.** Unlike prior works which applied a fixed *forget set* context in the prompt, or modified the query embeddings from original query to unlearn query, UNRE iteratively modifies query-related *forget subset* context embeddings from *forget piece* into *unlearn guiding context* (underlined). **Lower section: workflow of UNRE.** The circled step numbers show the UNRE modification progress, where step 2 to 5 happens entirely in the embedding space. Solid arrows represent communication between different components, and dashed arrows indicate state changes within the same component. The embedding model and LLM are off-the-shelf components. The legends at the top show where operations happen.

2024a)) to steer the LLM toward an unlearned output distribution (Wang et al., 2025b). Since the prompt encompasses all information provided to the model (Brown et al., 2020), inference-time unlearning, which is exemplified by input-based methods, operates during LLM inference with frozen weights and is generally regarded as suppression-intended unlearning (Ren et al., 2025). Several studies have also explored *query-adaptive* dynamic unlearning, for example, by leveraging a pre-trained reference model for real-time logit-difference computation (Ji et al., 2024), by applying inference-time prompt editing through pretrained rewrite agents (Sanyal & Mandal, 2025), or by modifying the decoding process (Deng et al., 2025; Gao et al., 2025). More recent work has introduced In-Context Unlearning (ICUL) (Pawelczyk et al., 2024; Takashiro et al., 2025), which highlights context manipulation as a new perspective within input-based unlearning, enabling preservation of LLM capabilities by retaining the original query, model weights, and architecture (Takashiro et al., 2025). **Meanwhile, unlearning requests will update frequently thus continual unlearning ability is important.** (Gao et al., 2025)

However, existing studies have not yet explored zero-shot LLM unlearning through contextual mechanisms, and particularly the query-adaptive dynamic contexts. This gap is important: in real-world dynamic settings such as privacy protection and bias mitigation, practitioners typically only have access to the *forget set*. When unlearning is required at inference time, lightweight methods that can rapidly adapt to changes in the *forget set* are especially valuable. Moreover, such approaches offer the additional advantage of not requiring any modification to the model parameters.

In this work, we present UNRE, a RAG-based method that refines query-retrieved contexts via gradient updates in the embedding space, and leverages these adapted contexts to guide the LLM toward generating outputs aligned with the unlearned distribution. The overview of the method is shown in Figure 1. Our main contributions are summarized below:

- We introduce UNRE, to the best of our knowledge, the first framework to leverage query-specific dynamic contexts for achieving zero-shot unlearning.
- We develop an online membership-inference-guided RAG architecture that first identifies the query-specific unlearning target, then optimizes the retrieved contexts for unlearning generation, thereby minimizing disruption to the LLM’s original capabilities.
- UNRE is 100% pre-processing free and query-adaptive, which better aligns with the dynamic requirements of real-world scenarios.

- Through extensive experiments on diverse unlearning tasks across multiple LLMs, UNRE demonstrates superior unlearning effectiveness, e.g., around 3 times stronger than fixed contextual unlearning, while largely preserving the model’s original performance by maintaining a similar perplexity score as the original model.

2 RELATED WORKS

Machine Unlearning Machine unlearning aims to remove specific behaviors or knowledge without broadly degrading model utility (Cao & Yang, 2015). Liu et al. (2024b) formulates the unlearning target as a confounder between an LLM’s input and output, and casts unlearning as a deconfounding process. Recent works have explored loss adjustment using only the *forget set*, thereby avoiding reliance on retain data or auxiliary agents (Wang et al., 2025b; Yao et al., 2024a). Zero-shot unlearning was introduced as a scenario where only the forget data is available (Foster et al., 2024). For instance, Gu et al. (2025) proposed to generate an anti-forget set to enhance fine-tuning-based unlearning. Other approaches, including PROD (Jiang et al., 2025), DPO (Rafailov et al., 2023), and NPO (Zhang et al., 2024), constrain unlearning with original model outputs to preserve overall performance. On the other hand, CEU (Entesari et al., 2025) flattens the forget set logits while enforcing a retain set performance lower bound to preserve model utility after tuning. In this work, we adopt a related loss-based formulation but operate solely on the *forget set*, eliminating the need for the retain data or auxiliary models.

Inference-time Unlearning Prompt engineering has emerged as a method for performing unlearning at inference time. For example, SPUL (Bhaila et al., 2025) trains soft prompts during an offline stage using a pre-established forget set and retain set, guiding LLMs to generate outputs that approximate a machine-unlearned distribution. *ECO* (Liu et al., 2024a) trained a classifier for unlearn-required prompt offline and a general corruption parameter that is applied to vectorized user input query in the embedding space to guide LLM to generate output in an unlearned distribution. Contrastive decoding methods, such as UCD and ULD (Suriyakumar et al., 2025; Ji et al., 2024), leverage logit differences between a small model trained on Forget and Retain Sets to guide unlearning. Since providing context adaptive to a specific query will make the LLM perform In-Context Learning (ICL) differently (Garg et al., 2022), and context can be updated by gradient in the embedding space during inference time (Zou et al., 2023), In-Context Unlearning (ICUL) (Pawelczyk et al., 2024) uses prompt context constructed from Forget and Retain Sets to prevent the generation of unwanted content. Vector steering techniques (Li et al., 2023; Rimsky et al., 2024; Arditì et al., 2024; Cao et al., 2024; Dunefsky & Cohan), including *InferAligner* (Wang et al., 2024a) and *FairSteer* (Li et al., 2025), inject offline-prepared steering vectors into LLM layers to influence inference. Other methods modify the decoder or employ multi-agent systems for unlearning (Wang et al., 2025a; Deng et al., 2025; Sanyal & Mandal, 2025). All these approaches, however, require either a retain set or offline training/tuning.

Retrieval Augmented Generation and Unlearning via RAG Retrieval-Augmented Generation (RAG) has seen significant advances in recent years, improving LLM performance by providing relevant external information during generation (Lewis et al., 2020b). A typical RAG pipeline involves chunking, embedding, retrieval, and reranking (Lewis et al., 2020a), and recent methods have focused on better aligning the retriever with the LLM. For example, *REPLUG* (Shi et al., 2024) tunes the retriever based on LLM’s likelihood signal to update the retrieval embeddings via gradient, which improves both perplexity and downstream accuracy. In this work, we adopt a multi-query RAG framework (Cheng et al., 2024) in which the embeddings are aligned with the LLM and receive gradient updates from its outputs, enabling more effective and adaptive retrieval during generation.

Several recent works have explored using RAG for unlearning at inference time. Wang et al. (2024b) constructs a retain set from the forget set offline and injects it into RAG for unlearning. *Eraser4RAG* (Wang et al., 2025d) trains a rewrite agent via reinforcement learning to transform retrieved forget data into retainable content, while *De-Indexing* (Vilella & Ruffo, 2025) reranks retrieved items to promote the retain set over the forget set. Similar to other prior works on inference-time unlearning, all these approaches require either a preprocessed retain set or prior agent training, limiting their applicability in scenarios where only the forget set is available.

3 METHODOLOGY

In this work, we propose an inference-time zero-shot unlearning framework UNRE based on RAG that operates solely with the *forget set*, requiring neither additional training nor fine-tuning throughout the workflow nor any architectural modifications to the LLM.

3.1 PROBLEM STATEMENT

The goal of query-adaptive zero-shot unlearning is to force the targeted LLM M to generate an output y_u in an unlearned token-sequence distribution, given only the *forget set* \mathbf{O} , where $\mathbf{O} = \{(\mathbf{o}_i)\}_{i=1}^n$ ($i \in [1, n]$) is the n chunk pieces among the *forget set* and \mathbf{o}_i is a piece in the *forget set*. Each forget piece $\mathbf{o}_i = \{\mathbf{x}_i, \mathbf{y}_i\}$ contains the *feature/example input* \mathbf{x} and *label* \mathbf{y} .

We consider a scenario where the UNRE owner has access to neither the *model training data* nor the *retain data* (in contrast to prior unlearning methods (Yao et al., 2024a)), and where the user query itself remains unaltered (Liu et al., 2024a)—with only the retrieved context being modified at inference time. To this end, the unlearning objective is to find a perturbed set $\tilde{\mathbf{O}}_q$, where \mathbf{O}_q is the query-related subset of \mathbf{O} , so that using $\tilde{\mathbf{O}}_q$ as context to constrain the LLM inference generation progress $y_u = M(q; \tilde{\mathbf{O}}_q)$.

3.2 METHOD OVERVIEW

We propose UNRE framework to find the proper $\tilde{\mathbf{O}}_q$. The overview of the method is shown in Figure 1. During the offline stage, \mathbf{O} will be input into the RAG, a technique that augments LLM generation through retrieving query-related contexts, stored in the RAG embedding vector database \mathbf{V}_R (steps 1 to 2 in Figure 1), which can be aligned with LLM embeddings \mathbf{V}_M (Cheng et al., 2024).

During inference, UNRE consists of the following stages. First, online membership inference (He et al., 2025; Fu et al., 2024) for \mathbf{O}_q (steps 1 to 2 in Figure 1). When receiving the user query q , the query will go through the LLM, generating a regular output y_q . Then, the input query q and output y_q will be sent to a RAG retrieval module, which will conduct a similarity search in \mathbf{V}_R (step 3 in Figure 1). If the similarity result is higher than a *threshold gate* τ , unlearning is required. Second, a dynamic unlearned context updating process for $\tilde{\mathbf{O}}_q$. This is achieved through gradient descent inside embedding space, aiming at optimizing the *unlearning loss function* (steps 4 to 5 in Figure 1) by using a *perturbation matrix* δ , which constrains the LLM output $y_u = M(q; \tilde{\mathbf{O}}_q)$ (step 6 in Figure 1) into an unlearned distribution. Note that ϕ is the original embedding model (such as e5).

3.3 PRE-CHECK GATE WITH ONLINE MEMBERSHIP INFERENCE

The pre-check progress aims to minimize the UNRE influence on the model’s original performance by shutting down the UNRE when not needed. We first obtain a regular output $y_q = M(q)$. By following a standard retrieval process in RAG (Lewis et al., 2020a), we compute its retrieval similarity to the *forget set* \mathbf{O} , as follows:

$$\max_{i \in [n]} \text{Similarity}(\phi(y_q), \phi(o_i)) < \tau. \quad (1)$$

The similarity threshold τ is a user-defined parameter, and ϕ is the embedding progress. We use *L2 distance in embedding* for similarity calculation. If the similarity is below τ , y_q is returned; otherwise, the UNRE pipeline starts.

3.4 UNLEARN PREPARATION OF UNRE

Thus, we have a real-time, query-specific *forget subset* \mathbf{O}_q through the retrieval progress,

$$\mathbf{O}_q = \{o_i \in \mathbf{O} : -d(\phi(y_q), \phi(o_i)) \geq \tau\}. \quad (2)$$

where d is the L2 distance between embedding vectors; and thus we have the embedding vector of sub-*forget set* $E_{\mathbf{O}} = \phi(\mathbf{O}_q)$. Through the retrieval similarity search during $\phi(y_q)$, we obtain the

Algorithm 1 Gradient-based Update in UNRE Embedding

Require: query q ; LLM \mathbf{M} ; matrixed forget context v_1 ; pre-set budget ε ; pre-set step size η ; gradient update steps $j \in [1, J]$; query-specific unlearn matrix δ ; token position $t \in [1, T]$ in a sentence; last-layer hidden state $h_\delta(t)$; last-layer logits $z_\delta(t)$

- 1: $\delta^{(0)} \leftarrow \mathbf{0}$; $y_q \leftarrow \mathbf{M}(q)$
- 2: *Get* $h_0(t)$, $z_0(t)$ while \mathbf{M} generating y_q
- 3: **for** $j = 1$ **to** J **do**
- 4: $\tilde{v}_1 \leftarrow v_1 + \delta^j$;
- 5: LLM output: $y_u \leftarrow \mathbf{M}(q \oplus \tilde{v}_1)$
- 6: *Get* $h_\delta(t)$, $z_\delta(t)$ while \mathbf{M} generating
- 7: Calculate $\bar{h}_\delta \leftarrow \frac{1}{T} \sum_{t=1}^T h_\delta(t)$;
- 8: Semantic Similarity: $S \leftarrow \text{Similarity}(\bar{h}_\delta, \phi(y_q))$
- 9: Calculate $\hat{z}_\delta(t)$ given $z_\delta(t)$; $\hat{z}_0(t)$ given $z_0(t)$
- 10: Next-Token Distribution Shift: $\mathcal{N} \leftarrow \frac{1}{T} \sum_{t=1}^T \cos(\hat{z}_\delta(t), \hat{z}_0(t))$
- 11: $\mathcal{L} \leftarrow \text{softplus}(\mathcal{N} - S)$
- 12: Gradient Update $\delta^{j+1} \leftarrow \text{PGD}(\mathcal{L})$
- 13: **end for**
- 14: **return** δ^J

initial RAG-retrieved, query-related embedding E_R ,

$$E_R = \text{TopK}(-d(\phi(y_q), \phi(o_i))), \quad (3)$$

where K is the RAG retrieval parameter defined by the user (Lewis et al., 2020b).

We obtain $E_1 = E_O \cup E_R$, and then form E_1 into matrix v_1 . Prior works like *PromptReps* and *HyDE* (Gao et al., 2023; Zhuang et al., 2024) have shown that looping back LLM output can enhance the retrieval progress. Starting from the retrieved query-related context and limiting the context example amount K can retain more of LLM’s original performance (Pawelczyk et al., 2024).

3.5 UPDATE EMBEDDING VECTORS TO UNLEARN IN UNRE

3.5.1 DESIGN OF LOSS FUNCTION FOR THE UNLEARNING OBJECTIVE

We optimize a perturbation matrix δ and feed the perturbed input $\tilde{v}_1 = v_1 + \delta$ into \mathbf{M} . The method is detailed in Algorithm 1. **In our setting, we perturb (applying δ to) the groundtruth y_i — the specific content that must be forgotten — following ICUL (Pawelczyk et al., 2024), which applies label reversal during its offline stage to enhance forgetting.**

At each output token position $t \in \{1, \dots, T\}$, the model outputs a last-layer hidden state $h_\delta(t)$ and logits $z_\delta(t)$ Yao et al. (2024a). **\oplus in line 5 represents the progress of projecting the embedding vector back to discrete tokens and combining it with the current text.** We start at $\delta = \mathbf{0}$, where we have $h_0(t)$, $z_0(t)$, as shown in lines 1 and 2 in Algorithm 1. The design of the UNRE loss is motivated by maintaining the *semantic meanings* while increasing *token distributional shift* (Sinha et al., 2025; Liu et al., 2024b; Wang et al., 2025c) of unlearned output y_u , as discussed below.

Sentence semantics For Semantic Similarity S , as illustrated in lines 7 to 8 of the Algorithm 1, we aggregate hidden states into a sentence vector and compute semantic similarity (higher is better).

Distributional shift of next-token predictions. As presented in lines 9 to 10 of Algorithm 1, we let $\hat{z}(t) = z(t)/\|z(t)\|_2$ denote the unit-direction of logits. We define \mathcal{N} , the expectation of the token-level directional discrepancy (lower is better). Consequently, we have the *loss function* (line 11 of Algorithm 1):

$$\mathcal{L}(\delta) = \text{softplus}(\mathcal{N} - S) = \log(1 + \exp(\mathcal{N} - S)) \quad (4)$$

Since loss adjustment can flexibly realize diverse unlearning objectives (Wang et al., 2025b), we generalize the loss function \mathcal{L} to support a broad range of unlearning tasks (e.g., copyright, privacy) by tuning task-specific parameters and integrating the ECO loss formulation (Liu et al., 2024a), as detailed in Appendix E.1.2.

3.5.2 GRADIENT-BASED UPDATE IN CONTEXT EMBEDDING

We employ Projected Gradient Descent (PGD) (Madry et al., 2018) to update gradients in the embedding space, while constraining the update region to avoid the *forget set* embedding $e = \phi(O)_q$. Specifically, we optimize the perturbation δ (line 10 in Algorithm 1) to minimize the loss \mathcal{L} .

$$\delta^{(j+1)} = \Pi_{\{\delta: \|\delta\|_2 \leq \varepsilon, \min_{e \in E} \|v_1 + \delta - e\|_2 \geq \tau\}} \left(\delta^{(j)} - \eta \nabla_{\delta} \mathcal{L}(\delta^{(j)}) \right), \quad (5)$$

where $\delta^{(j)}$ has the same dimension as v_1 , denotes the query-specific unlearning perturbation matrix at PGD step j ; $\nabla_{\delta} \mathcal{L}(\delta^{(j)})$ is the gradient of the loss \mathcal{L} evaluated at $\delta^{(j)}$; η is the learning rate for the gradient update; Π is the projection operator onto the perturbation ball with budget ε specified by the UNRE owner; j is the iteration index; and J is the total number of PGD steps.

3.6 UNRE UNLEARNING INFERENCE

After obtaining δ^J from Algorithm 1, we construct the perturbed matrix $v_c = v_1 + \delta^J$ for final inference. UNRE then constrains the LLM using updated contexts $C = \tilde{\mathbf{O}}_q$ decoded from v_c , thereby guiding the model to generate unlearned outputs $y_u = \mathbf{M}(q; C)$.

Finally, the pre-check procedure described in Section 3.3 is started again to determine whether another run of UNRE is necessary, ensuring that the final LLM output y_u exhibits no similarity to the *forget set* \mathbf{O} .

4 EXPERIMENT

Overview In this section, we evaluate UNRE across a range of unlearning tasks, including *Entity Unlearning* and *Copyright Content Unlearning*, using the *TOFU* (Maini et al., 2024), *RW KU* (Jin et al., 2024), and *HP* (Eldan & Russinovich, 2023) datasets. We further assess its performance on *context unlearning* under varying context lengths, comparing against a state-of-the-art in-context unlearning method (Pawelczyk et al., 2024). Additional tasks and results are provided in Appendix. All experiments are conducted on Nvidia L40S GPUs.

Baseline Methods We compare UNRE against a diverse set of unlearning baselines, grouped into three categories. *Gradient-based methods* include Gradient Ascent (**GA**) (Maini et al., 2024), Grad-Diff (**GD**) (Maini et al., 2024), KL minimization (**KL**) (Maini et al., 2024), Large Language Model Unlearning (**LLMU**) (Yao et al., 2024a), and **Mismatch** (Yao et al., 2024b), as well as regularized GA variants **GAGDR** and **GAKLR** (Shi et al., 2025). *Preference-based methods* include Preference Optimization (**PO**) (Maini et al., 2024), Direct Preference Optimization (**DPO**) (Maini et al., 2024), Negative Preference Optimization (**NPO**) (Zhang et al., 2024), and the regularized NPO variants **NPOGDR** and **NPOKLR** (Shi et al., 2025), together with the forget-only loss-adjustment method **FLAT** (Wang et al., 2025b). *Tuning-free methods* include In-Context Unlearning (**ICUL**) (Pawelczyk et al., 2024), **ECO** (Liu et al., 2024a), **GUARD** (Deng et al., 2025), and **Prompt/Output-Filtering** strategies (Deng et al., 2025; Pawelczyk et al., 2024). We include more baselines and their descriptions in Appendix.

4.1 ENTITY UNLEARNING

4.1.1 TOFU 1% SPLIT

We evaluate entity unlearning on the **TOFU 1% Split** benchmark (Maini et al., 2024) Following prior work, we first fine-tune each base LLMs on the full TOFU training set to obtain the *Original LLM*; the *Retained LLM* is fine-tuned on the split, which serves as the reference model. We report the 1% forget split and use LLMs of Falcon3-7B, Llama3.2-3B and Qwen2.5-7B, as summarized in Table 1.

Metrics We adopt the official TOFU evaluation metrics. **Forget Quality (FQ)** is defined as the p -value from a Kolmogorov–Smirnov test applied to the Truth Ratio distributions of the unlearned and retained models on the forget set; higher values indicate stronger unlearning performance. **Model Utility (MU)** is computed as the harmonic mean of Answer Probability, Truth Ratio, and ROUGE-L

Table 1: **TOFU 1% split**. Performance of our method and baseline methods on the TOFU dataset using three base LLMs (Falcon3-7B, Llama3.2-3B and Qwen2.5-7B). FQ, MU, F-RL, and R-RL denote *forget quality*, *model utility*, *ROUGE-L on the forget set*, and *ROUGE-L on the retain set*, respectively. We include the Original LLM and the Retained LLM (trained on retain set) for reference.

Method	Falcon3-7B-Instruct				Llama3.2-3B-Instruct				Qwen2.5-7B-Instruct			
	FQ↑	MU↑	F-RL↓	R-RL↑	FQ↑	MU↑	F-RL↓	R-RL↑	FQ↑	MU↑	F-RL↓	R-RL↑
Original LLM	0.0067	0.6644	0.8612	0.8030	0.0067	0.5752	0.9913	0.9778	0.0067	0.6054	0.9719	0.9219
Retained LLM	1.0	0.6647	0.3792	0.7998	1.0	0.6018	0.4088	0.9866	1.0	0.5910	0.3794	0.8958
GA	0.0067	0.6663	0.7379	0.8041	0.0067	0.5754	0.8112	0.9735	0.0541	0.5887	0.4723	0.8837
KL	0.0067	0.6653	0.7347	0.7943	0.0066	0.5759	0.8331	0.9755	0.0970	0.5876	0.4613	0.8820
GD	0.0286	0.6535	0.7058	0.8195	0.0066	0.5747	0.8359	0.9771	0.0286	0.5929	0.4745	0.8848
LLMU	0.0287	0.6544	0.7589	0.8183	0.0143	0.5680	0.9913	0.9765	0.0286	0.5656	0.4774	0.5823
PO	0.0067	0.6625	0.8290	0.8084	0.0143	0.5678	0.9913	0.9774	0.0067	0.6152	0.7387	0.8459
DPO	0.0286	0.6535	0.7058	0.8195	0.0065	0.5766	0.7379	0.9769	0.0067	0.5766	0.7379	0.5259
NPO	0.0067	0.6656	0.7432	0.7958	0.0067	0.5768	0.7866	0.9765	0.0143	0.5539	0.4055	0.5258
FLAT	0.0030	0.6659	0.7013	0.7994	0.0066	0.5766	0.7379	0.9769	0.0286	0.5971	0.5079	0.9032
ICUL	0.0286	0.6641	0.4059	0.8028	0.0143	0.5751	0.5614	0.9778	0.0143	0.6054	0.4539	0.9217
Prompt	0.0970	0.6644	0.4045	0.8030	0.0143	0.5753	0.8635	0.9777	0.0067	0.6053	0.5552	0.9218
GUARD	0.0541	0.6643	0.3115	0.8029	0.5786	0.5752	0.3764	0.9776	0.2656	0.6052	0.3691	0.9219
UnRe (Ours)	0.0611	0.6644	0.2824	0.8030	0.6012	0.5752	0.3298	0.9778	0.2977	0.6054	0.3169	0.9219

across the subsets retain, real authors, world facts, where higher scores reflect better utility preservation. We also report **F-RL** (ROUGE-L on the forget set; lower is better) and **R-RL** (ROUGE-L on the retain set; higher is better).

Results It can be seen that UNRE demonstrates strong unlearning performance while preserving model utility across modern LLMs. On Llama3.2-3B and Qwen2.5-7B, it achieves model utility (MU) scores of 0.5752 and 0.6054, staying within 0.28% and 1.6% of the best train-time baselines (Original/ICUL/GUARD). At the same time, UNRE attains superior Forget Quality (FQ), reaching 0.6012 on Llama3.2-3B and 0.2977 on Qwen2.5-7B, surpassing GUARD, while remaining competitive on Falcon3-7B (0.0611) and outperforming gradient-based baselines. Across all models, it maintains a favorable forget–retain trade-off, achieving the lowest forget–retain loss (F-RL) while keeping retain–retain loss (R-RL) at the Original level. Overall, UNRE provides effective unlearning with minimal impact on model utility across different LLMs.

4.1.2 REAL-WORLD KNOWLEDGE UNLEARNING (RWKU)

We also evaluate entity unlearning on the **RWKU** benchmark (Jin et al., 2024), as a *test-only* suite. The *Original LLM* (Before) denotes the base model without unlearning, and our method is applied at inference time using a forget-only retrieval corpus derived from RWKU materials. We report results on LLaMA-3-8B-Instruct and LLaMA-3.1-8B-Instruct, as presented in Table 2.

Metrics We adopt the official RWKU metrics. **Forget** reports ROUGE-L on Fill-in-the-Blank and QA probes over the forget targets (FB/QA; lower is better); **AA** denotes adversarial probes in robustness analyses. **Neighbor** is ROUGE-L on probes about entities adjacent to the forget targets and reflects locality (higher is better). **MIA** reports membership inference on forget- and retain-like samples via **FM** (higher is better) and **RM** (lower is better). **Utility** measures general capabilities on reasoning, truthfulness, factual QA, and fluency (**Rea, Tru, Fac, Flu**; higher is better).

Results Across both LLMs, UNRE delivers the strongest *forgetting* while preserving *locality*, *privacy*, and *utility*. On LLaMA-3-8B, it reduces **Forget–QA** to 39.8, outperforming most baselines, and increases **Neighbor–QA** to 78.1 (vs. 76.5 for GAGDR, the base baseline), indicating reduced collateral forgetting. For **MIA**, UNRE achieves higher **FM** (268.7, above NPO as the best baseline) and lower **RM**, reflecting weaker membership signals on the forget set and fewer false positives on *retain-like data*. **Utility** is maintained showing minimal degradation to general model capabilities. On LLaMA-3.1-8B, UNRE further lowers **Forget–AA** to 38.7, and improves **Neighbor–FB** to 74.0 (surpassing best baseline NPO_{GDR}). For **MIA**, it achieves higher **FM** and the best **RM**, reflecting effective unlearning without falsely flagging *retain data*. Model **Utility** remains robust for LLaMA-3.1-8B as well, with overall trends comparable to baseline performance.

Table 2: **RWKU**. We report *Forget* (FB/QA/AA/All, ↓), *Neighbor* (FB/QA/All, ↑), *MIA* (FM↑/RM↓), and *Utility* (Rea/Tru/Fac/Flu, ↑).

(a) LLaMA-3-8B-Instruct

	Forget ↓			Neighbor ↑		MIA		Utility ↑			
Method	FB	QA	AA	FB	QA	FM ↑	RM ↓	Rea	Tru	Fac	Flu
Before	85.6	70.3	74.7	93.1	82.0	236.5	230.9	41.0	36.4	53.7	704.6
GA	72.0	64.6	68.5	85.0	74.7	241.4	234.6	40.4	37.6	49.6	710.3
GAGDR	72.6	64.0	69.7	86.2	76.5	242.8	236.8	39.6	36.8	50.4	710.3
GAKLR	70.7	57.5	69.9	80.5	70.5	242.4	230.8	41.5	35.6	54.0	704.4
NPO	46.6	39.0	35.3	79.2	70.9	263.3	241.4	40.5	36.0	56.7	695.9
NPOGDR	52.2	43.9	42.9	82.5	70.5	254.5	240.1	39.6	37.2	51.4	708.2
NPOKLR	52.5	40.6	43.2	83.2	72.1	253.0	236.9	40.9	35.4	54.2	704.9
UnRe (Ours)	44.8	39.8	34.9	88.4	78.1	267.7	236.2	40.6	36.0	53.7	704.6

(b) LLaMA-3.1-8B-Instruct

	Forget ↓			Neighbor ↑		MIA		Utility ↑			
Method	FB	QA	AA	FB	QA	FM ↑	RM ↓	Rea	Tru	Fac	Flu
Before	63.9	65.1	69.5	74.1	69.8	223.5	218.2	42.2	35.4	61.2	695.2
GA	50.7	45.4	61.2	45.6	37.2	248.9	241.9	43.2	35.8	48.7	726.6
GAGDR	55.4	49.6	63.9	60.2	53.5	239.8	231.3	44.2	35.0	53.9	718.5
GAKLR	62.7	49.9	66.4	67.9	61.2	235.8	223.0	42.6	35.4	59.0	682.1
NPO	35.7	40.2	39.0	67.3	66.2	241.4	220.5	42.5	35.6	61.8	684.2
NPOGDR	42.4	37.2	42.0	74.0	66.7	236.3	220.1	43.0	35.4	60.8	698.8
NPOKLR	40.6	41.4	42.2	73.3	69.9	234.4	218.8	42.3	35.4	61.5	695.1
UnRe (Ours)	39.2	37.9	38.7	74.0	68.6	242.4	220.1	42.3	35.4	61.2	695.0

4.2 COPYRIGHTED CONTENT UNLEARNING

We use **Harry Potter and the Sorcerer’s Stone** (Eldan & Russinovich, 2023) (HP) as copyrighted content to be forgotten, constructing **forget** and **retain** splits by extracting 400 chunks from the book for the *forget set* and sampling 400 paragraphs from C4 for the *retain set*. The LLM is fine-tuned on the *forget set* to simulate memorization, while the original pretrained checkpoint serves as the retained baseline.

Metrics We report the **Forget Quality Gap (FQ Gap)** defined over BLEU and ROUGE-L differences between the unlearned and the retained model on *the split forget set*, together with **Perplexity (PPL)** (Jelinek et al., 1977) and the average zero-shot accuracy (**Avg. Acc.**) across nine standard tasks as a model-utility proxy. We evaluate on OPT-2.7B and Llama2-7B models for better comparison with prior works.

Results It can be seen from the results that UNRE achieves effective unlearning without compromising model utility in general. In the HP setting, it consistently enforces strong forgetting while preserving general capabilities. Operating entirely at inference time, the framework activates conservatively only on copyright-relevant queries, ensuring that LLM generation text quality (PPL) and zero-shot accuracy remain aligned with the original checkpoint across architectures. This demonstrates the core goal of inference-time unlearning: *eliminate targeted knowledge while maintaining unrelated model capabilities*.

Besides, it can be observed that prior methods, which do not explicitly balance forgetting and utility, typically fail in one of two ways: (i) improving the forget score but degrading fluency or accuracy, or (ii) preserving general performance while leaving residual memorization. By contrast, UNRE successfully preserves the retained model’s utility profile while removing reproduction of the copyrighted text. Prompt- or filter-based baselines largely leave non-trigger inputs unchanged and fail to provide targeted suppression, whereas optimization-based methods can achieve forgetting but often at the expense of the generation text quality. The results highlight that, as a lightweight, training-free method, UNRE effectively performs copyright-content unlearning in HP, achieving targeted knowledge removal while retaining overall generation performance, outperforming prior methods.

Table 3: **HP unlearning** on OPT-2.7B and Llama2-7B. Lower FQ Gap/PPL and higher Avg. Acc. are better.

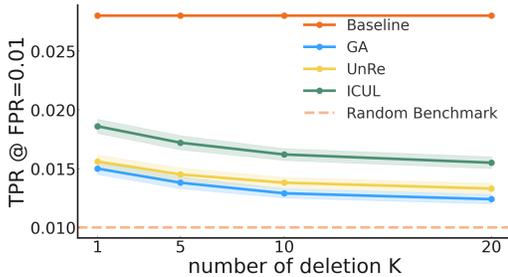
Method	OPT-2.7B			Llama2-7B		
	FQ Gap ↓	PPL ↓	Avg. Acc. ↑	FQ Gap ↓	PPL ↓	Avg. Acc. ↑
Original LLM	1.5346	15.6314	0.4762	3.6594	8.9524	0.5617
Retained LLM	0.0000	14.3190	0.4686	0.0000	8.7070	0.5599
KL	2.7301	16.1592	0.4688	0.4225	9.4336	0.5509
GD	2.3439	16.1972	0.4690	0.5304	9.1797	0.4902
Mismatch	1.4042	15.7507	0.4679	0.4647	8.9906	0.5593
LLMU	2.4639	15.8398	0.4656	0.1985	9.0530	0.5503
PO	2.1601	14.8960	0.4583	0.5124	8.8364	0.5532
DPO	2.2152	16.8396	0.4621	0.2924	8.9597	0.5614
NPO	1.2611	19.6637	0.4644	0.5151	9.0397	0.5609
FLAT	1.4089	15.5543	0.4686	0.2265	8.9906	0.5580
ICUL	1.0121	15.6314	0.4762	2.5585	8.9524	0.5617
GUARD	0.6314	15.6314	0.4762	0.1367	8.9524	0.5617
UnRe (Ours)	0.6112 ± 0.0011	15.6314	0.4762	0.1207 ± 0.0008	8.9524	0.5617

Method	Accuracy ↑	TPR @ FPR=0.01 ↓
Baseline	90.4%	0.0267
ICUL	89.8%	0.0179
UnRe	90.0 ± 0.2%	0.0147

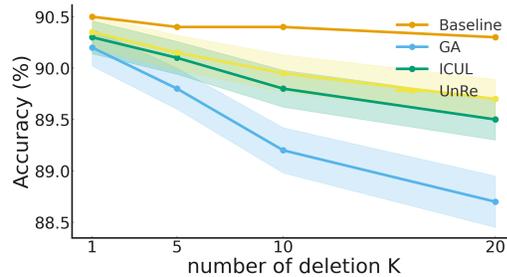
(a) Unlearning results at $K = 10$.

Method	Accuracy ↑	TPR @ FPR=0.01 ↓
Baseline	90.4%	0.0267
ICUL	90.2%	0.0183
UnRe	90.3 ± 0.2%	0.0153

(b) Unlearning results at $K = 5$



(c) Unlearning performance (TPR).



(d) Accuracy.

Figure 2: Evaluate unlearning for different numbers of deletion requests (1, 5, 10, 20).

4.3 CONTEXTUAL UNLEARNING COMPARISON FOR DIFFERENT CONTEXT LENGTHS

We follow the ICUL (Pawelczyk et al., 2024) setup and adopt its *LiRA-Forget* protocol (Carlini et al., 2022) to quantify unlearning. We evaluate inference-time unlearning across varying context lengths, considering 5 and 10 deletions (i.e., $K = 5$ and $K = 10$ retrieved context examples). As shown in Figure 2.

Metrics The **TPR @ FPR=0.01** measures the true positive rate of a likelihood-ratio test distinguishing an unlearned model from a retained trained model on the forget points (*lower is better*). **Accuracy** reflects standard test performance on held-out data, serving as a utility indicator (*higher is better*). Effective unlearning is indicated by TPR values approaching the benchmark while maintaining accuracy close to the baseline. **Baseline** refers to the original fine-tuned model without any unlearning.

Results UNRE demonstrates the intended behavior of inference-time unlearning using in-context examples: it largely preserves task accuracy, consistently outperforming ICUL across varying context lengths and approaching the performance of GA. While ICUL’s forgetting improves with longer contexts, its overall unlearning effectiveness remains substantially below that of UNRE.

5 CONCLUSION

In this work, we propose UNRE, a novel retrieval-based framework for dynamic, query-adaptive zero-shot unlearning in LLMs. Unlike prior approaches that rely on fixed prompts or static context injection, UNRE leverages query-adaptive dynamic contexts to achieve inference-time unlearning without any offline preparation. The framework first employs online membership inference to guide retrieval from the *forget set*, adapting context to each query, and then applies gradient-based perturbations to the retrieved embeddings to steer the LLM’s outputs toward an unlearned distribution. Empirical results across multiple LLMs and unlearning tasks demonstrate that UNRE effectively removes targeted knowledge while preserving the model’s original capabilities. Notably, it operates without pretraining or retain sets, making it particularly suitable for lightweight, real-world unlearning scenarios where the *forget set* is frequently updated. Overall, UNRE illustrates that dynamic context can enable efficient, query-adaptive zero-shot unlearning during LLM inference.

ETHICS STATEMENT

We adhere to the ICLR Code of Ethics. UNRE is an inference-time, training-free unlearning controller that operates *only* with the forget set and leaves the base model’s parameters unchanged; a conservative pre-check gate prevents activation on benign inputs. As a result, the method targets removal/suppression of copyrighted passages and hazardous knowledge while preserving general utility, thereby *reducing* potential harm rather than introducing new risks. Our experiments rely on standard public benchmarks (e.g., Harry Potter excerpts for copyright unlearning; WMDP for hazardous-knowledge attenuation) and do not involve human subjects or the collection of personal data; no copyrighted material is redistributed. We release code and prompts with safeguards aimed at preventing misuse (e.g., documentation on intended use and limitations). Overall, UNRE is designed to strengthen ethical deployment by enabling targeted forgetting without degrading unrelated capabilities.

REPRODUCIBILITY STATEMENT

All experimental settings (datasets, splits, preprocessing, model variants, hyperparameters, training schedules, and evaluation protocols) are described in detail in Section 4. We conduct all experiments on a single node equipped with $4 \times$ NVIDIA L40S GPUs. We submit the code in the supplementary material, which includes a fully specified runtime environment and scripts to reproduce results.

REFERENCES

- 540
541
542 Sk Miraj Ahmed, Umit Yigit Basaran, Dripta S. Raychaudhuri, Arindam Dutta, Ro-
543 hit Kundu, Fahim Faisal Niloy, Basak Guler, and Amit K. Roy-Chowdhury. Towards
544 source-free machine unlearning. In *IEEE/CVF Conference on Computer Vision and Pat-
545 tern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pp. 4948–4957.
546 Computer Vision Foundation / IEEE, 2025. doi: 10.1109/CVPR52734.2025.00466.
547 URL [https://openaccess.thecvf.com/content/CVPR2025/html/Ahmed_](https://openaccess.thecvf.com/content/CVPR2025/html/Ahmed_Towards_Source-Free_Machine_Unlearning_CVPR_2025_paper.html)
548 [Towards_Source-Free_Machine_Unlearning_CVPR_2025_paper.html](https://openaccess.thecvf.com/content/CVPR2025/html/Ahmed_Towards_Source-Free_Machine_Unlearning_CVPR_2025_paper.html).
- 549
550 Andy Ardit, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel
551 Nanda. Refusal in language models is mediated by a single direction. In Amir Globersons, Lester
552 Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang
553 (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on Neural
554 Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 -
555 15, 2024*, 2024. URL [http://papers.nips.cc/paper_files/paper/2024/hash/](http://papers.nips.cc/paper_files/paper/2024/hash/f545448535dfde4f9786555403ab7c49-Abstract-Conference.html)
556 [f545448535dfde4f9786555403ab7c49-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2024/hash/f545448535dfde4f9786555403ab7c49-Abstract-Conference.html).
- 557
558 Clark Barrett, Brad Boyd, Elie Bursztein, Nicholas Carlini, Brad Chen, Jihye Choi, Amrita Roy
559 Chowdhury, Mihai Christodorescu, Anupam Datta, Soheil Feizi, et al. Identifying and mitigating
560 the security risks of generative ai. *Foundations and Trends® in Privacy and Security*, 6(1):1–52,
561 2023.
- 562
563 Karuna Bhaila, Minh-Hao Van, and Xintao Wu. Soft prompting for unlearning in large language
564 models. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Proceedings of the 2025 Conference
565 of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human
566 Language Technologies, NAACL 2025 - Volume 1: Long Papers, Albuquerque, New Mexico, USA,
567 April 29 - May 4, 2025*, pp. 4046–4056. Association for Computational Linguistics, 2025. doi:
568 10.18653/V1/2025.NAACL-LONG.204. URL [https://doi.org/10.18653/v1/2025.](https://doi.org/10.18653/v1/2025.naacl-long.204)
569 [naacl-long.204](https://doi.org/10.18653/v1/2025.naacl-long.204).
- 570
571 Lucas Bourtole, Varun Chandrasekaran, Christopher A. Choquette-Choo, Hengrui Jia, Adelin
572 Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *2021 IEEE
573 Symposium on Security and Privacy (SP)*, pp. 141–159, 2021. doi: 10.1109/SP40001.2021.00019.
- 574
575 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-
576 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-
577 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,
578 Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler,
579 Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCand-
580 lish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot
581 learners. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan,
582 and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual
583 Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12,
584 2020, virtual*, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html)
585 [1457c0d6bfc4967418bfb8ac142f64a-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bfb8ac142f64a-Abstract.html).
- 586
587 Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In *2015
588 IEEE Symposium on Security and Privacy, SP 2015, San Jose, CA, USA, May 17-21, 2015*, pp.
589 463–480. IEEE Computer Society, 2015. doi: 10.1109/SP.2015.35. URL [https://doi.org/](https://doi.org/10.1109/SP.2015.35)
590 [10.1109/SP.2015.35](https://doi.org/10.1109/SP.2015.35).
- 591
592 Yuanpu Cao, Tianrong Zhang, Bochuan Cao, Ziyi Yin, Lu Lin, Fenglong Ma, and Jinghui
593 Chen. Personalized steering of large language models: Versatile steering vectors through
594 bi-directional preference optimization. In Amir Globersons, Lester Mackey, Danielle Bel-
595 grave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances
596 in Neural Information Processing Systems 38: Annual Conference on Neural Informa-
597 tion Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15,
598 2024*, 2024. URL [http://papers.nips.cc/paper_files/paper/2024/hash/](http://papers.nips.cc/paper_files/paper/2024/hash/58cbe393b4254da8966780a40d023c0b-Abstract-Conference.html)
599 [58cbe393b4254da8966780a40d023c0b-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2024/hash/58cbe393b4254da8966780a40d023c0b-Abstract-Conference.html).

- 594 Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramèr.
595 Membership inference attacks from first principles. In *43rd IEEE Symposium on Security and*
596 *Privacy, SP 2022, San Francisco, CA, USA, May 22-26, 2022*, pp. 1897–1914. IEEE, 2022. doi:
597 10.1109/SP46214.2022.9833649. URL [https://doi.org/10.1109/SP46214.2022.](https://doi.org/10.1109/SP46214.2022.9833649)
598 9833649.
- 599 Huiqiang Chen, Tianqing Zhu, Xin Yu, and Wanlei Zhou. Zero-shot machine unlearning with
600 proxy adversarial data generation. In *Proceedings of the Thirty-Fourth International Joint Con-*
601 *ference on Artificial Intelligence, IJCAI 2025, Montreal, Canada, August 16-22, 2025*, pp. 339–
602 347. ijcai.org, 2025. doi: 10.24963/IJCAI.2025/39. URL [https://doi.org/10.24963/](https://doi.org/10.24963/ijcai.2025/39)
603 [ijcai.2025/39](https://doi.org/10.24963/ijcai.2025/39).
- 604 Nuo Chen, Yan Wang, Haiyun Jiang, Deng Cai, Yuhan Li, Ziyang Chen, Longyue Wang, and
605 Jia Li. Large language models meet harry potter: A dataset for aligning dialogue agents with
606 characters. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Associa-*
607 *tion for Computational Linguistics: EMNLP 2023*, pp. 8506–8520, Singapore, December 2023.
608 Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.570. URL
609 <https://aclanthology.org/2023.findings-emnlp.570/>.
- 610 Xin Cheng, Xun Wang, Xingxing Zhang, Tao Ge, Si-Qing Chen, Furu Wei, Huishuai Zhang,
611 and Dongyan Zhao. xrag: Extreme context compression for retrieval-augmented gen-
612 eration with one token. In Amir Globersons, Lester Mackey, Danielle Belgrave, An-
613 gela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances in*
614 *Neural Information Processing Systems 38: Annual Conference on Neural Information*
615 *Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15,*
616 *2024, 2024*. URL [http://papers.nips.cc/paper_files/paper/2024/hash/](http://papers.nips.cc/paper_files/paper/2024/hash/c5cf13bfd3762821ef7607e63ee90075-Abstract-Conference.html)
617 [c5cf13bfd3762821ef7607e63ee90075-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2024/hash/c5cf13bfd3762821ef7607e63ee90075-Abstract-Conference.html).
- 618 Vikram S. Chundawat, Ayush K. Tarun, Murari Mandal, and Mohan S. Kankanhalli. Zero-shot
619 machine unlearning. *IEEE Trans. Inf. Forensics Secur.*, 18:2345–2354, 2023. doi: 10.1109/TIFS.
620 2023.3265506. URL <https://doi.org/10.1109/TIFS.2023.3265506>.
- 621 Zhijie Deng, Chris Yuhao Liu, Zirui Pang, Xinlei He, Lei Feng, Qi Xuan, Zhaowei Zhu, and Jiaheng
622 Wei. Guard: Generation-time llm unlearning via adaptive restriction and detection. *arXiv preprint*
623 *arXiv:2505.13312*, 2025.
- 624 Guangyao Dou, Zheyuan Liu, Qing Lyu, Kaize Ding, and Eric Wong. Avoiding copyright infringe-
625 ment via large language model unlearning. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.),
626 *Findings of the Association for Computational Linguistics: NAACL 2025, Albuquerque, New*
627 *Mexico, USA, April 29 - May 4, 2025*, pp. 5176–5200. Association for Computational Linguis-
628 tics, 2025. doi: 10.18653/v1/2025.FINDINGS-NAACL.288. URL [https://doi.org/10.](https://doi.org/10.18653/v1/2025.findings-naacl.288)
629 [18653/v1/2025.findings-naacl.288](https://doi.org/10.18653/v1/2025.findings-naacl.288).
- 630 Jacob Dunefsky and Arman Cohan. One-shot optimized steering vectors mediate safety-relevant
631 behaviors in llms. In *Second Conference on Language Modeling*.
- 632 Javid Ebrahimi, Anyi Rao, Daniel Lowd, and Dejing Dou. Hotflip: White-box adversarial examples
633 for text classification. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th An-*
634 *nual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia,*
635 *July 15-20, 2018, Volume 2: Short Papers*, pp. 31–36. Association for Computational Linguistics,
636 2018. doi: 10.18653/v1/P18-2006. URL <https://aclanthology.org/P18-2006/>.
- 637 Ronen Eldan and Mark Russinovich. Who’s harry potter? approximate unlearning in llms. *CoRR*,
638 [abs/2310.02238](https://arxiv.org/abs/2310.02238), 2023. doi: 10.48550/ARXIV.2310.02238. URL [https://doi.org/10.](https://doi.org/10.48550/arXiv.2310.02238)
639 [48550/arXiv.2310.02238](https://doi.org/10.48550/arXiv.2310.02238).
- 640 Taha Entesari, Arman Hatami, Rinat Khaziev, Anil Ramakrishna, and Mahyar Fazlyab. Con-
641 strained entropic unlearning: A primal-dual framework for large language models. *arXiv preprint*
642 *arXiv:2506.05314*, 2025.
- 643 Jack Foster, Kyle Fogarty, Stefan Schoepf, Cengiz Öztireli, and Alexandra Brintrup. Zero-shot
644 machine unlearning at scale via lipschitz regularization. *CoRR*, [abs/2402.01401](https://arxiv.org/abs/2402.01401), 2024. doi: 10.
645 [48550/ARXIV.2402.01401](https://arxiv.org/abs/2402.01401). URL <https://doi.org/10.48550/arXiv.2402.01401>.

- 702 15, 2024, 2024. URL [http://papers.nips.cc/paper_files/paper/2024/](http://papers.nips.cc/paper_files/paper/2024/hash/blf78dfc9ca0156498241012aec4efa0-Abstract-Datasets_and_)
703 hash/blf78dfc9ca0156498241012aec4efa0-Abstract-Datasets_and_
704 Benchmarks_Track.html.
- 705
706 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman
707 Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel,
708 and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In
709 Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-
710 Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Confer-
711 ence on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020,*
712 *virtual*, 2020a. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html)
713 6b493230205f780e1bc26945df7481e5-Abstract.html.
- 714 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,
715 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe
716 Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural
717 Information Processing Systems*, volume 33, pp. 9459–9474, 2020b.
- 718 Kenneth Li, Oam Patel, Fernanda B. Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-
719 time intervention: Eliciting truthful answers from a language model. In Alice Oh, Tris-
720 tan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Ad-
721 vances in Neural Information Processing Systems 36: Annual Conference on Neural Infor-
722 mation Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,*
723 *2023*, 2023. URL [http://papers.nips.cc/paper_files/paper/2023/hash/](http://papers.nips.cc/paper_files/paper/2023/hash/81b8390039b7302c909cb769f8b6cd93-Abstract-Conference.html)
724 81b8390039b7302c909cb769f8b6cd93-Abstract-Conference.html.
- 725 Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D.
726 Li, Ann-Kathrin Dombrowski, Shashwat Goel, Gabriel Mukobi, Nathan Helm-Burger, Rassin
727 Lababidi, Lennart Justen, Andrew Bo Liu, Michael Chen, Isabelle Barrass, Oliver Zhang, Xi-
728 aoyuan Zhu, Rishub Tamirisa, Bhruu Bharathi, Ariel Herbert-Voss, Cort B Breuer, Andy
729 Zou, Mantas Mazeika, Zifan Wang, Palash Oswal, Weiran Lin, Adam Alfred Hunt, Justin
730 Tienken-Harder, Kevin Y. Shih, Kemper Talley, John Guan, Ian Steneker, David Campbell, Brad
731 Jokubaitis, Steven Basart, Stephen Fitz, Ponnurangam Kumaraguru, Kallol Krishna Karmakar,
732 Uday Tupakula, Vijay Varadharajan, Yan Shoshitaishvili, Jimmy Ba, Kevin M. Esvelt, Alexandr
733 Wang, and Dan Hendrycks. The WMDP benchmark: Measuring and reducing malicious use
734 with unlearning. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria
735 Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International
736 Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*,
737 pp. 28525–28550. PMLR, 21–27 Jul 2024. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v235/li24bc.html)
738 v235/li24bc.html.
- 739 Yichen Li, Zhiting Fan, Ruizhe Chen, Xiaotang Gai, Luqi Gong, Yan Zhang, and Zuozhu Liu.
740 Fairsteer: Inference time debiasing for llms with dynamic activation steering. In Wanxiang
741 Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Findings of
742 the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August
743 1, 2025*, pp. 11293–11312. Association for Computational Linguistics, 2025. URL [https://](https://aclanthology.org/2025.findings-acl.589/)
744 aclanthology.org/2025.findings-acl.589/.
- 745 Chris Yuhao Liu, Yaxuan Wang, Jeffrey Flanigan, and Yang Liu. Large language model un-
746 learning via embedding-corrupted prompts. In Amir Globersons, Lester Mackey, Danielle
747 Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Ad-
748 vances in Neural Information Processing Systems 38: Annual Conference on Neural Infor-
749 mation Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15,*
750 *2024*, 2024a. URL [http://papers.nips.cc/paper_files/paper/2024/hash/](http://papers.nips.cc/paper_files/paper/2024/hash/d6359156e0e30b1caa116a4306b12688-Abstract-Conference.html)
751 d6359156e0e30b1caa116a4306b12688-Abstract-Conference.html.
- 752 Yujian Liu, Yang Zhang, Tommi S. Jaakkola, and Shiyu Chang. Revisiting who’s harry pot-
753 ter: Towards targeted unlearning from a causal intervention perspective. In Yaser Al-Onaizan,
754 Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical
755 Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pp. 8708–8731. Association for Computational Linguistics, 2024b. doi:

- 756 10.18653/V1/2024.EMNLP-MAIN.495. URL <https://doi.org/10.18653/v1/2024.>
757 emnlp-main.495.
758
- 759 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
760 Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*,
761 2017.
- 762 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
763 Towards deep learning models resistant to adversarial attacks. In *International Conference on*
764 *Learning Representations*, 2018.
- 765
- 766 Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C Lipton, and J Zico Kolter. Tofu: A task
767 of fictitious unlearning for llms. *arXiv preprint arXiv:2401.06121*, 2024.
- 768
- 769 Fabio Motoki, Valdemar Pinho Neto, and Victor Rodrigues. More human than human: measuring
770 chatgpt political bias. *Public Choice*, 198(1):3–23, 2024.
- 771
- 772 Martin Pawelczyk, Seth Neel, and Himabindu Lakkaraju. In-context unlearning: Language models
773 as few-shot unlearners. In *Forty-first International Conference on Machine Learning, ICML 2024,*
774 *Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL [https://openreview.](https://openreview.net/forum?id=GKcwl8XC9)
775 [net/forum?id=GKcwl8XC9](https://openreview.net/forum?id=GKcwl8XC9).
- 776
- 777 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Ermon, and
778 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model.
779 In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine
780 (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural*
781 *Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -*
782 *16, 2023*, 2023. URL [http://papers.nips.cc/paper_files/paper/2023/hash/](http://papers.nips.cc/paper_files/paper/2023/hash/a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html)
783 [a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html).
- 784
- 785 Jie Ren, Yue Xing, Yingqian Cui, Charu C. Aggarwal, and Hui Liu. Sok: Machine unlearning for
786 large language models. *CoRR*, abs/2506.09227, 2025. doi: 10.48550/ARXIV.2506.09227. URL
787 <https://doi.org/10.48550/arXiv.2506.09227>.
- 788
- 789 Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner.
790 Steering llama 2 via contrastive activation addition. In Lun-Wei Ku, Andre Martins, and Vivek
791 Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational*
792 *Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*,
793 pp. 15504–15522. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.
794 [acl-long.828](https://doi.org/10.18653/v1/2024.acl-long.828). URL <https://doi.org/10.18653/v1/2024.acl-long.828>.
- 795
- 796 Debdeep Sanyal and Murari Mandal. Agents are all you need for llm unlearning. In *Second Confer-*
797 *ence on Language Modeling*, 2025.
- 798
- 799 Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke
800 Zettlemoyer, and Wen-tau Yih. REPLUG: retrieval-augmented black-box language models.
801 In Kevin Duh, Helena Gómez-Adorno, and Steven Bethard (eds.), *Proceedings of the 2024*
802 *Conference of the North American Chapter of the Association for Computational Linguistics:*
803 *Human Language Technologies (Volume 1: Long Papers), NAACL 2024, Mexico City, Mex-*
804 *ico, June 16-21, 2024*, pp. 8371–8384. Association for Computational Linguistics, 2024. doi:
805 [10.18653/v1/2024.naacl-long.463](https://doi.org/10.18653/v1/2024.naacl-long.463). URL [https://doi.org/10.18653/v1/2024.](https://doi.org/10.18653/v1/2024.naacl-long.463)
806 [naacl-long.463](https://doi.org/10.18653/v1/2024.naacl-long.463).
- 807
- 808 Weijia Shi, Jaechan Lee, Yangsibo Huang, Sadhika Malladi, Jieyu Zhao, Ari Holtzman, Daogao
809 Liu, Luke Zettlemoyer, Noah A. Smith, and Chiyuan Zhang. MUSE: machine unlearning six-
way evaluation for language models. In *The Thirteenth International Conference on Learn-*
ing Representations, ICLR 2025, Singapore, April 24-28, 2025. OpenReview.net, 2025. URL
<https://openreview.net/forum?id=TArMA033BU>.
- 809
- 808 Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference at-
809 tacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy (SP)*,
pp. 3–18, 2017. doi: 10.1109/SP.2017.41.

- 810 Yash Sinha, Murari Mandal, and Mohan S. Kankanhalli. Unstar: Unlearning with self-taught
811 anti-sample reasoning for llms. *Trans. Mach. Learn. Res.*, 2025, 2025. URL [https://](https://openreview.net/forum?id=mNXCvIKZbI)
812 openreview.net/forum?id=mNXCvIKZbI.
813
- 814 Minkyoo Song, Hanna Kim, Jaehan Kim, Seungwon Shin, and Soeul Son. Refusal is not an option:
815 Unlearning safety alignment of large language models. In *34th USENIX Security Symposium*
816 (*USENIX Security 25*), pp. 319–338, 2025.
- 817 Vinith M Suriyakumar, Ayush Sekhari, and Ashia Wilson. Ucd: Unlearning in llms via contrastive
818 decoding. *arXiv preprint arXiv:2506.12097*, 2025.
- 819 Shota Takashiro, Takeshi Kojima, Andrew Gambardella, Qi Cao, Yusuke Iwasawa, and Yutaka
820 Matsuo. Answer when needed, forget when not: Language models pretend to forget via in-
821 context knowledge unlearning. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and
822 Mohammad Taher Pilehvar (eds.), *Findings of the Association for Computational Linguistics:*
823 *ACL 2025*, pp. 24872–24885, Vienna, Austria, July 2025. Association for Computational Lin-
824 guistics. ISBN 979-8-89176-256-5. doi: 10.18653/v1/2025.findings-acl.1276. URL [https://](https://aclanthology.org/2025.findings-acl.1276/)
825 aclanthology.org/2025.findings-acl.1276/.
826
- 827 Salvatore Vilella and Giancarlo Ruffo. (de)-indexing and the right to be forgotten. *CoRR*,
828 [abs/2501.03989](https://doi.org/10.48550/ARXIV.2501.03989), 2025. doi: 10.48550/ARXIV.2501.03989. URL [https://doi.org/10.](https://doi.org/10.48550/ARXIV.2501.03989)
829 [48550/arXiv.2501.03989](https://doi.org/10.48550/ARXIV.2501.03989).
- 830 Haoran Wang, Xiong Xiao Xu, Baixiang Huang, and Kai Shu. Privacy-aware decoding: Mit-
831 igating privacy leakage of large language models in retrieval-augmented generation. *CoRR*,
832 [abs/2508.03098](https://doi.org/10.48550/ARXIV.2508.03098), 2025a. doi: 10.48550/ARXIV.2508.03098. URL [https://doi.org/10.](https://doi.org/10.48550/ARXIV.2508.03098)
833 [48550/arXiv.2508.03098](https://doi.org/10.48550/ARXIV.2508.03098).
- 834 Pengyu Wang, Dong Zhang, Linyang Li, Chenkun Tan, Xinghao Wang, Mozhi Zhang, Ke Ren,
835 Botian Jiang, and Xipeng Qiu. Inferaligner: Inference-time alignment for harmlessness through
836 cross-model guidance. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceed-*
837 *ings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP*
838 *2024, Miami, FL, USA, November 12-16, 2024*, pp. 10460–10479. Association for Computa-
839 tional Linguistics, 2024a. doi: 10.18653/v1/2024.EMNLP-MAIN.585. URL [https://doi.](https://doi.org/10.18653/v1/2024.emnlp-main.585)
840 [org/10.18653/v1/2024.emnlp-main.585](https://doi.org/10.18653/v1/2024.emnlp-main.585).
- 841 Shang Wang, Tianqing Zhu, Dayong Ye, and Wanlei Zhou. When machine unlearning meets
842 retrieval-augmented generation (RAG): keep secret or forget knowledge? *CoRR*, [abs/2410.15267](https://doi.org/10.48550/ARXIV.2410.15267),
843 2024b. doi: 10.48550/ARXIV.2410.15267. URL [https://doi.org/10.48550/ARXIV.](https://doi.org/10.48550/ARXIV.2410.15267)
844 [2410.15267](https://doi.org/10.48550/ARXIV.2410.15267).
- 845 Yaxuan Wang, Jiaheng Wei, Chris Yuhao Liu, Jinlong Pang, Quan Liu, Ankit Shah, Yujia Bao, Yang
846 Liu, and Wei Wei. LLM unlearning via loss adjustment with only forget data. In *The Thirteenth*
847 *International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.*
848 *OpenReview.net*, 2025b. URL <https://openreview.net/forum?id=6ESRicalFE>.
849
- 850 Yaxuan Wang, Jiaheng Wei, Chris Yuhao Liu, Jinlong Pang, Quan Liu, Ankit Parag Shah, Yujia Bao,
851 Yang Liu, and Wei Wei. Llm unlearning via loss adjustment with only forget data. In *Proceedings*
852 *of the International Conference on Learning Representations*, 2025c.
- 853 Yujing Wang, Hainan Zhang, Liang Pang, Yongxin Tong, Binghui Guo, Hongwei Zheng, and Zhim-
854 ing Zheng. Learning to erase private knowledge from multi-documents for retrieval-augmented
855 large language models. *CoRR*, [abs/2504.09910](https://doi.org/10.48550/ARXIV.2504.09910), 2025d. doi: 10.48550/ARXIV.2504.09910. URL
856 <https://doi.org/10.48550/ARXIV.2504.09910>.
857
- 858 Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. In Amir Globersons,
859 Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng
860 Zhang (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on*
861 *Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December*
862 *10 - 15, 2024*, 2024a. URL [http://papers.nips.cc/paper_files/paper/2024/](http://papers.nips.cc/paper_files/paper/2024/hash/be52acf6bccf4a8c0a90fe2f5cfcead3-Abstract-Conference.html)
863 [hash/be52acf6bccf4a8c0a90fe2f5cfcead3-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2024/hash/be52acf6bccf4a8c0a90fe2f5cfcead3-Abstract-Conference.html).
- Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. In *NeurIPS*, 2024b.

864 Charles Yu, Sullam Jeong, Anish Kasi, Pengfei Yu, and Heng Ji. Unlearning bias in language
865 models by partitioning gradients. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki
866 (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 6032–6048,
867 Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.
868 findings-acl.375. URL <https://aclanthology.org/2023.findings-acl.375/>.

869 Xiaojian Yuan, Tianyu Pang, Chao Du, Kejiang Chen, Weiming Zhang, and Min Lin. A closer look
870 at machine unlearning for large language models. In *The Thirteenth International Conference*
871 *on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025.
872 URL <https://openreview.net/forum?id=Q1MHvGmhyT>.

873 Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. Negative preference optimization: From catastro-
874 phic collapse to effective unlearning. *CoRR*, abs/2404.05868, 2024. doi: 10.48550/ARXIV.
875 2404.05868. URL <https://doi.org/10.48550/arXiv.2404.05868>.

876 Shengyao Zhuang, Xueguang Ma, Bevan Koopman, Jimmy Lin, and Guido Zuccon. Promptreps:
877 Prompting large language models to generate dense and sparse representations for zero-shot doc-
878 ument retrieval. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings*
879 *of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024,*
880 *Miami, FL, USA, November 12-16, 2024*, pp. 4375–4391. Association for Computational Lin-
881 guistics, 2024. doi: 10.18653/v1/2024.EMNLP-MAIN.250. URL [https://doi.org/10.](https://doi.org/10.18653/v1/2024.emnlp-main.250)
882 [18653/v1/2024.emnlp-main.250](https://doi.org/10.18653/v1/2024.emnlp-main.250).

883 Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson.
884 Universal and transferable adversarial attacks on aligned language models. *arXiv preprint*
885 *arXiv:2307.15043*, 2023.
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

918 APPENDIX

919
920 A SUMMARY OF APPENDIX

921 We include the following supplementary materials that expand on our methods, experimental setups,
922 and evaluations.

923 **B LLM Usage Disclosure** - We detailed how we used LLM during the conduct of this project.

924 **C Hyperparameter** - We show hyperparameters we used in the experiments.

925 **D Extension Literature Reviews and Preliminaries** - We provided more details on the loss
926 function and preliminaries.

927 **E Additional Details of Methodology** - We provide the details of method scalability under dif-
928 ferent tasks, [alignment method](#), and [comparison with prior RAG-based unlearning](#), as well as
929 illustrating the inference workflow for better understanding.

930 **F Additional Experiment Settings** - We provide more explanations of the experiment settings.

931 **G Additional Experiments** - We provide a detailed comparison of different models (OPT fam-
932 ily and LLaMA family, as well as Mistral) with different tasks and datasets, as well as time
933 complexity, to show the effectiveness of our methods under different scenarios.

934 **H Hyperparameter Sensitivity Analysis** - We provide analysis of different hyperparameter sen-
935 sitivity.

936 **I Ablation Study** - We conduct ablation studies of different hyperparameter selections.

937 **J Visualization** - We visualize the dynamic contexts for better understanding the effect of our
938 methods.

939
940
941
942
943
944 B LLM USAGE DISCLOSURE

945 Large Language Models (LLMs) were used to aid in the writing and polishing of the manuscript.
946 Specifically, we used an LLM to assist in refining the language, improving readability, and ensuring
947 clarity in various sections of the paper. The model helped with tasks such as sentence rephrasing,
948 grammar checking, and enhancing the overall flow of the text.

949 It is important to note that the LLM was not involved in the ideation, research methodology, or
950 experimental design. All research concepts, ideas, and analyses were developed and conducted by
951 the authors. The contributions of the LLM were solely focused on improving the linguistic quality
952 of the paper, with no involvement in the scientific content or data analysis.

953 The authors take full responsibility for the content of the manuscript, including any text generated
954 or polished by the LLM. We have ensured that the LLM-generated text adheres to ethical guidelines
955 and does not contribute to plagiarism or scientific misconduct.

956
957
958
959 C HYPERPARAMETERS

960 We list all hyperparameters *used in this paper* in Table 4.

- 961 • **PGD step size** η : step length for each gradient update in the embedding-space PGD. We
962 treat η as a tunable hyperparameter.
- 963 • **PGD iteration count** J : total number of projected updates. Tuned for the utility/forgetting
964 trade-off.
- 965 • **PGD radius** ϵ : ℓ_2 budget that bounds the perturbation norm; enforced by projection.
- 966 • **Distance lower bound** τ_{dist} : optional constraint that keeps the (poisoned) context represen-
967 tation at least τ_{dist} away from the forget-set representation during projection.
- 968 • **Top-K retrieved contexts** K : number of passages retrieved per query in RAG; we expose
969 K as a user-level RAG knob and tune it for each dataset.

- **Gate threshold** τ : pre-check membership/similarity threshold that decides whether to *activate* our unlearning correction; larger values trigger more aggressively.

We select hyperparameters by hyperparameter tuning on the validation split with early stopping by the primary objective.

Table 4: Hyperparameters

Symbol	Name	Role	Status / Our setting
η	PGD step size	gradient update step in PGD	Hyperparameter Tuning
J	PGD iterations	# projected updates	HyperParameter Tuning
ϵ	PGD radius	ℓ_2 budget for projection	HyperParameter Tuning
τ_{dist}	distance lower bound	projection constraint	[0.0, 1.0]
K	Top-K contexts	RAG retrieval count	{3, 5, 10}
τ	gate threshold	trigger for applying UNRE	[0.0, 1.0]

We set these hyperparameters based on prior works by Lewis et al. (2020b); Madry et al. (2017). The settings can be easily adjusted according to the practical requirements.

D EXTENSION LITERATURE REVIEWS AND PRELIMINARIES

D.1 LOSS FUNCTION OF UNLEARNING WORKS

Design of Loss Our loss $\mathcal{L}(\delta)$ couples (i) semantic preservation and (ii) distributional shift of next-token predictions. It is inspired by inference-time steering ideas (e.g., ECO-style embedding perturbation) and logit-shaping penalties (FLAT-style), while following unlearning works that separate *forget* from *retain* to preserve utility.

Relation to prior unlearning and steering. ECO performs inference-time corruption in embedding space without updating weights; FLAT-style methods penalise logit geometry; loss-adjustment unlearning enforces pressure on the forget set while regularising retain performance. Our \mathcal{L} inherits the inference-time, weight-frozen setting, but acts on the retrieved context and explicitly couples semantic preservation with logit-direction shift.

Reference loss templates (for citation). We summarize the canonical objectives we are inspired from; each is shown in a compact form.

$$\mathcal{L}_{\text{UL}}^f(\theta) = \frac{1}{T} \sum_{t=1}^T -\log\left(1 - p_{\theta}(y_t^* | y_{<t}, x)\right), \quad (6)$$

$$\mathcal{L}_{\text{KL}}^r(\theta) = \frac{1}{T} \sum_{t=1}^T \text{KL}\left(p_{\theta}(\cdot | y_{<t}, x) \parallel p_{\text{ref}}(\cdot | y_{<t}, x)\right), \quad (7)$$

$$\mathcal{L}_{\text{ECO}}(\delta) = \frac{1}{T} \sum_{t=1}^T \text{KL}\left(\text{softmax}(z_{\delta}(t)/\tau) \parallel \text{softmax}(z_{\text{ret}}(t)/\tau)\right), \quad (8)$$

$$\mathcal{L}_{\text{COS}}(\delta) = \frac{1}{T} \sum_{t=1}^T \cos(\widehat{z}_{\delta}(t), \widehat{z}_0(t)), \quad \widehat{z}(\cdot) = \frac{z(\cdot)}{\|z(\cdot)\|_2}, \quad (9)$$

$$\mathcal{L}_{\text{KL}}^{\text{shift}}(\delta) = \frac{1}{T} \sum_{t=1}^T \text{KL}\left(\text{softmax}(z_{\delta}(t)/\tau) \parallel \text{softmax}(z_0(t)/\tau)\right). \quad (10)$$

$$\mathcal{L}_{\text{SEM}}(\delta) = 1 - \text{Similarity}(\bar{h}_\delta, \phi(y_0)), \quad (11)$$

Mapping to our loss. In Equation 4, the term \mathcal{N} instantiates a logit-shift penalty (e.g., Equation 9 or Equation 10), while S is the complement of Equation 11; optional retain regularization Equation 7 can be added if needed.

D.2 PRELIMINARY

Unlearning objective extension UNRE maintains sentence-level semantics while pushing away next-token directional predictions from those of y_q :

$$\bar{h}_\delta := \frac{1}{T} \sum_{t=1}^T h_\delta(t), \quad S := \text{sim}(\bar{h}_\delta, \varphi(y_q)), \quad \pi := \frac{1}{T} \sum_{t=1}^T \cos\left(\frac{z_\delta(t)}{\|z_\delta(t)\|_2}, \frac{z_0(t)}{\|z_0(t)\|_2}\right), \quad (12)$$

$$\mathcal{L}(\delta) = \text{softplus}(\pi - S), \quad \delta^{(j+1)} = \Pi_{\|\delta\|_2 \leq \varepsilon}(\delta^{(j)} - \eta \nabla_\delta \mathcal{L}(\delta^{(j)})). \quad (13)$$

Final delivery. After J steps, $v_c := v_1 + \delta^{(J)}$ is either (i) delivered directly to V_M (continuous injection), or (ii) Forward if previously decoded to text and concatenated to the prompt.

E ADDITIONAL DETAILS OF METHODOLOGY

E.1 LOGIT NORMALIZATION AND CENTERING

For improved invariance and stability one may replace $z(t)$ by a normalized direction:

$$\textbf{Unit-only: } \hat{z}(t) = \frac{z(t)}{\|z(t)\|_2}. \quad (14)$$

Time pooling We pool per-step directions and compare only the pooled vectors:

$$\hat{\tilde{z}}_\delta = \frac{1}{T_\delta} \sum_{t=1}^{T_\delta} \hat{z}_\delta(t), \quad \hat{\tilde{z}}_0 = \frac{1}{T_0} \sum_{t=1}^{T_0} \hat{z}_0(t), \quad R_{\text{dist}}^{\text{pool}} = \cos(\hat{\tilde{z}}_\delta, \hat{\tilde{z}}_0).$$

Use $R_{\text{dist}}^{\text{pool}}$ in Equation 4 as a drop-in replacement for R_{dist} .

E.1.1 EXTENSION OF LOSS

$$\mathcal{L}(\delta) = \log(1 + \exp(\mathcal{N} - S)) \quad (15)$$

$$\mathcal{L}(\delta) = \max(\mathcal{N} - S, 0) + \log(1 + \exp(-|\mathcal{N} - S|)) \quad (16)$$

E.1.2 ADAPT LOSS TO DIFFERENT UNLEARNING TASKS

Following the loss design of Liu et al. (2024a), the UNRE Loss \mathcal{L} can be extended into:

$$\mathcal{L}_{\text{unified}}(\delta) = \underbrace{\text{softplus}(\mathcal{N} - S)}_{\text{UNRE base}} + \alpha \cdot \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} \rho_a \ell_a(s_a(\tilde{y}), \omega_a) \quad (17)$$

where

- $\alpha \geq 0$ is a mixing coefficient that weights the ECO-style target term against the UNRE base term $\text{softplus}(\mathcal{N} - S)$.
- \mathcal{A} is the index set of task assessments (e.g., multiple-choice accuracy, BERTScore, ROUGE, ASG, PII hit-rate, etc.); $a \in \mathcal{A}$ indexes one assessment.

Algorithm 2 UNREembedding alignment

Require: query q ; LLM M with Embedding Model ϕ_M ; LLM Embedding Space V_M ; matrixed forget context v_1 ; optimize matrix δ^J ;

- 1: **if** EmbeddingAlignment = *False* **then**
- 2: $C \leftarrow \text{Decode}(v_1)$
- 3: $v_1^M \leftarrow \phi_M(C)$
- 4: **else**
- 5: $v_1^M \leftarrow v_1$
- 6: **end if**
- 7: $\tilde{v}_1 \leftarrow v_1^M + \delta^J$
- 8: Optimize δ^J
- 9: UNREed LLM output: $y_u \leftarrow M_{V_M}(\phi_M(q) + \tilde{v}_1)$
- 10: *Restore* UNREed Context: C ; UNREed Embedding Matrix: \tilde{v}_1
- 11: **return** UNREed LLM output: y_u

- $s_a(\tilde{y})$ denotes the score of the assessment a computed on the generated sequence \tilde{y}
- ω_a is the desired target for assessment a that represents a *retain-like* state in the sense of ECO (Liu et al., 2024a).
- $\rho_a \geq 0$ is an optional weight for assessment a .
- $\ell_a(\cdot, \cdot)$ is a penalty for assessment a

E.1.3 TASK-SPECIFIC INSTANTIATION

We keep the loss form in Equation 17 unchanged and only plug in task-specific assessments $\langle \mathcal{A}, s_a, \omega_a, \ell_a, \rho_a \rangle$.

E.1.4 CONTEXT FORMING IN DIFFERENT SCENARIOS

Since there is an option that RAG can send the embedding vectors to LLM embeddings (Cheng et al., 2024), UNRE can align with the LLM so that it can directly connect to the embedding space.

Based on the scenarios (when the LLM embeddings can not be aligned), UNREowner can decide if they want to deliver \tilde{v}_1 to LLM Embedding V_M or align V_M with V_R , showing in line 2 in Algorithm 2. Line 8 can be referred to Alg. 1. Meanwhile, we store \tilde{v}_1 and C , linked to v_1 , as a UNRE-ed set that can be transferred to other unlearning methods for use as the *retain set*.

E.2 JUSTIFICATIONS

We view the frozen LLM M as a differentiable map from the perturbed context embedding $v_1 + \delta$ to the sequence of next-token logits $\{z_\delta(t)\}_{t=1}^T$ defined in Eq. (12). Thus, the UNRE loss in Eq. (4) can be written explicitly as

$$L(\delta) = L(M(q; v_1 + \delta)),$$

Consequently, the update in Eq. 5 performs gradient descent on this output-level unlearning objective with respect to δ via the chain rule, while keeping weights fixed, in the same idea as ECO’s embedding-corrupted prompts and soft-prompt unlearning (Liu et al., 2024a; Bhaila et al., 2025).

E.3 ADDITIONAL REVIEWS AND NOVELTY COMPARISON

E.3.1 COMPARISON WITH PRIOR RAG-BASED UNLEARNINGS

Prerequisites for Unlearning As previously discussed, unlearning typically requires two datasets. The forget set contains knowledge that the LLM has memorized during training but is now required to forget; the retain set contains knowledge that the LLM should continue to remember and output correctly. In this paper, we consider the scenario of the zero-shot continual unlearning setting, in which only the forget set is available during unlearning.

RAG-based Unlearning via Offline Reverse-Generated Retain Set (Wang et al., 2024b). In our formulation, let \mathbf{O} be the *forget set*. Wang et al. (2024b) constructs a *retain set* \mathcal{R} from \mathbf{O} offline (through reverse generation), and then injects it into the RAG pipeline

Formally, the offline stage applies a reverse-generation transformation

$$\mathcal{R} = T_{\text{offline}}(\mathbf{O}), \quad (18)$$

where T_{offline} is implemented by using an LLM agent and prompt engineering. Once this procedure finishes, the offline reverse-generated *retain set* \mathcal{R} is frozen (fixed).

At inference time, given a query q , Wang et al. (2024b) uses a standard query-dependent retriever on the fixed corpus \mathcal{R} .

$$\mathcal{R}_q = \text{Retrieval}(q; \mathcal{R}), \quad (19)$$

$$y_u^{\text{prior}} = \mathbf{M}(q; \mathcal{R}_q). \quad (20)$$

where \mathcal{R}_q denotes the retrieved retain context produced by the standard RAG retrieval process (Eq. 19) under the method of Wang et al. (2024b). It is frozen/fixed during the online stage because \mathcal{R} is processed during the offline stage and is frozen during the online stage. Thus, the retrieved context \mathcal{R}_q cannot adapt content to the specific query, since \mathcal{R} does not consider real-time q . The method cannot further refine the content of the context.

Prior RAG-based methods vs. UNRE. In contrast, UNRE dynamically optimizes the content (tokens) of a query-specific perturbed context set $\tilde{\mathbf{O}}_q$ at inference time. Let \mathbf{O} denote the *forget set*; \mathbf{O}_q be the query-related subset retrieved at inference time (similar to the standard retrieval, see Eq. 2 in the Methodology Section), UNRE then performs an online constrained optimization towards the content of \mathbf{O}_q into $\tilde{\mathbf{O}}_q$:

$$\tilde{\mathbf{O}}_q = \arg \min_{\tilde{\mathbf{O}}_q} \mathcal{L}(\mathbf{M}(q; \tilde{\mathbf{O}}_q), \mathbf{O}_q) \quad \text{s.t.} \quad \|\tilde{\mathbf{O}}_q - \mathbf{O}_q\| \leq \epsilon, \quad (21)$$

where \mathcal{L} is our unlearning loss (Eq. 4) and the gradient optimization (Eq. 5) is carried out in the embedding space. The final LLM unlearned output is

$$y_u^{\text{UNRE}} = \mathbf{M}(q; \tilde{\mathbf{O}}_q). \quad (22)$$

Unlike prior RAG-based unlearning, the unlearned context $\tilde{\mathbf{O}}_q$ (compared to \mathcal{R}_q , which is frozen during online) is not a fixed offline object but a *query-adaptive* solution computed at inference time and directly targeted at the desired LLM unlearned output.

In this case, *forget set* \mathbf{O} can be updated frequently based on the task requirements, and no offline preparation process is required. Meanwhile, the retrieved contexts are more directly targeting the inference-time unlearning goal.

E.4 ALIGNMENT DETAILS

Following xRAG (Cheng et al., 2024), for each document (forget piece) o_i , the dense retrieval feature is

$$e_R(o_i) = \phi(o_i) \in \mathbf{V}_R. \quad (23)$$

xRAG introduces a *modality projector* with trainable weight θ

$$B_\theta : \mathbf{V}_R \rightarrow \mathbf{V}_M, \quad (24)$$

which is the only trainable component for projection; both ϕ and \mathbf{M} remain frozen. The projector maps $e_R(o_i)$ into the LLM embedding space as a single “document token”:

$$e_M(o_i) = B_\theta(e_R(o_i)) \in \mathbf{V}_M. \quad (25)$$

Given a query q with token embeddings $\phi_M(q) \in \mathbf{V}_M$, xRAG feeds the LLM with the input embedding sequence $\phi_M(q) \oplus e_M(o_i)$.

Consequently, \mathbf{M} treats $e_M(o_i)$ indistinguishably from a standard document token in its native representation space.

The projector B_θ is trained (with ϕ and \mathbf{M} frozen) by combining a language-modeling loss and a self-distillation loss, as detailed in Eq. (2)–(3) of (Cheng et al., 2024).

1188 F ADDITIONAL EXPERIMENT SETTINGS

1189

1190 **More introductions on baselines.**

1191

1192 **GA:** gradient-ascent on forget data to suppress target likelihood.

1193

1194 **KL:** GA with KL-to-reference regularization to preserve utility.

1195

1196 **GD:** gradient-based unlearning with direct loss on forget and a retain-side utility term (lightweight GA variant)

1197

1198 **LLMU:** a train-time unlearning recipe combining GA, random-mismatch loss, and KL-to-original for stability

1199

1200 **PO:** preference-style optimization that downranks forget-consistent responses relative to retain-consistent ones

1201

1202 **DPO:** direct preference optimization adapted to unlearning (no reward model)

1203

1204 **NPO:** negative preference optimization to avoid GA collapse and improve the forget/utility trade-off

1205

1206 **FLAT:** forget-data-only loss adjustment (no retain data / no reference model)

1207

1208 **ICUL:** in-context unlearning via specially constructed contexts and a likelihood-ratio signal at inference time

1209

1210 **Prompt:** rule-based *output filtering / guardrails* that refuse/deflect on forget-related queries

1211

1212 **GUARD:** detection + adaptive restriction during decoding to block forbidden tokens/semantics

1213

1214 **ECO:** embedding-corrupted prompts gated by a prompt classifier to enforce an “unlearned state” at inference

1215

1216 **GAGDR/GAKLR:** GA augmented with (i) gradient-direction regularization (GDR) or (ii) KL-to-retained anchoring (KLR) to stabilize utility (regularized GA variants).

1217

1218 **NPOGDR/NPOKLR:** NPO with the same (GDR/KLR) retain-side regularizers (regularized NPO variants).

1219

1220 **Mismatch:** context-mismatch baseline pairing queries with intentionally mismatched passages to reduce recall of copyrighted/entity text

1221

1222 G ADDITIONAL EXPERIMENT RESULTS

1223

1224 G.1 MUSE-NEWS UNLEARNING

1225

1226 **Experiment Setup** We evaluate on **MUSE-News** with its two tasks: *VerbMem* (verbatim memorization) and *KnowMem* (knowledge memorization).

1227

1228 It can be seen that in the MUSE benchmark (Table 5), UNRE also outperforms most baseline methods.

1229

1230 G.2 KNOWLEDGE UNLEARNING WMDP

1231

1232 Table 6 presents the unlearning results of GRUN and provides a direct comparison against the other evaluated methods.

1233

1234 G.3 TIME COMPLEXITY

1235

1236 UNRE performs only a single lightweight offline step: it embeds the forget set and builds a small index used by the pre-check gate. During inference, it simply conducts a gated similarity check, and all model weights remain unchanged. In contrast, baselines such as LLMU require fine-tuning, and ICUL requires reverse-generation to construct context data.

1237

1238 We report the time complexity for the experiments of Table 3 in Table 7. Compared with ICUL Pawelczyk et al. (2024), UNRE increases the per-query runtime by only about 4–5 percent (1.31

1240

1242 Table 5: MUSE-News results (official four metrics). Lower is better for VerbMem/KnowMem on
 1243 D_f (forget set), higher is better for KnowMem on D_r (retain set), and PrivLeak should be close to
 1244 0.

1246	Method	VerbMem on D_f ↓	KnowMem on D_f ↓	KnowMem on D_r ↑	PrivLeak
1247	Original LLM	58.4	63.9	55.2	-99.8
1248	Retained LLM	20.8	33.1	55.0	0.0
1249	Task Vectors	56.3	63.7	54.6	-99.8
1250	WHP	19.7	21.2	28.3	109.6
1251	GA	0.0	0.0	0.0	17.0
1252	GD	4.9	27.5	6.7	109.4
1253	KL	27.4	50.2	44.8	-96.1
1254	NPO	0.0	0.0	0.0	15.0
1255	NPO-RT	1.2	54.6	40.5	105.8
1256	Mismatch	42.8	52.6	45.7	-99.8
1257	FLAT (TV)	1.7	13.6	31.8	45.4
1258	FLAT (KL)	0.0	0.0	0.0	58.9
1259	FLAT (JS)	1.9	36.2	38.5	47.1
1260	FLAT (Pearson)	1.6	0.0	0.2	26.8
1261	ICUL	10.7	19.7	55.2	-99.8
1262	Output Filtering	1.1	0.3	55.2	-99.8
1263	Prompt	15.4	47.9	55.2	-99.6
1264	GUARD	4.3	4.9	55.2	109.6
1265	UnRe	4.0	33.2	55.2	-99.8

1264 Table 6: WMDP results reported by GRUN. Bio/Cyber are accuracies (0–1).
 1265

1266	Model	Method	Bio ↓	Cyber ↓	MMLU ↑
1267	Llama 3.1	Before	0.696	0.418	0.611
1268		Vanilla	0.494	0.337	0.581
1269		GRUN	0.372	0.293	0.577
1270	Mistral v0.1	Before	0.668	0.437	0.581
1271		Vanilla	0.256	0.252	0.529
1272		GRUN	0.293	0.278	0.535

1273 vs. 1.25). However, because UNRE requires no offline stage, its overall time for one epoch is
 1274 significantly lower. Compared to LLMU (Yao et al., 2024a), which relies heavily on an offline
 1275 stage, UnRe achieves roughly a 44 percent reduction in total runtime.
 1276

1277 Table 7: Time Complexity (seconds, averaged)
 1278

1281	Method	Offline Total	Online Total	Overall Runtime for One Epoch
1282	LLMU	1684	493	2177
1283	ICUL	317	534	851
1284	UnRe	0	637	637

1285 G.4 UNLEARNING PERFORMANCE COMPARING WITH RAG

1286 Since the traditional RAG-based unlearning methods can just reduce the RAG augmenting perform-
 1287 ance through reranking or unlearning through loading retain document (Wang et al., 2024b), or
 1288 require a *retain set*, thus UNRE is not comparable with traditional RAG-based unlearning methods.
 1289

1290 G.5 TOFU 1% SPLIT ON MORE MODELS

1291 **UNRE preserves model utility.** As shown in Table 8, UNRE incurs almost no degradation in model
 1292 utility compared to the original/retained references. On **Llama2-7B**, UnRe attains a top-2 MU, on
 1293

Table 8: TOFU 1% split. Performance of our method and baseline methods on the TOFU dataset using two base LLMs (Llama2-7B and Phi-1.5B). FQ, MU, F-RL, and R-RL denote *forget quality*, *model utility*, *ROUGE-L on the forget set*, and *ROUGE-L on the retain set*, respectively.

Method	Llama2-7B				Phi-1.5B			
	FQ↑	MU↑	F-RL↓	R-RL↑	FQ↑	MU↑	F-RL↓	R-RL↑
Original LLM	4.4883e-06	0.6346	0.9851	0.9833	0.0013	0.5184	0.9607	0.9199
Retained LLM	1.0	0.6267	0.4080	0.9833	1.0	0.5233	0.4272	0.9269
GA	0.0143	0.6333	0.4862	0.9008	0.0013	0.5069	0.5114	0.8048
KL	0.0068	0.6300	0.5281	0.9398	0.0030	0.5047	0.5059	0.8109
GradDiff	0.0068	0.6320	0.4773	0.8912	0.0030	0.5110	0.4996	0.8496
PO	0.0541	0.6308	0.3640	0.8811	0.0286	0.5127	0.3170	0.7468
Mismatch	0.0143	0.6304	0.9406	0.9741	0.0030	0.5225	0.9612	0.9194
LLMU	0.0030	0.5999	0.4891	0.9236	0.0143	0.5083	0.3380	0.7685
ICUL	0.0005	0.6239	0.4772	0.9818	0.0286	0.5195	0.0564	0.9276
Output Filtering	0.0002	0.6239	0.0	0.9818	0.00002	0.5195	0.0	0.9276
Prompt	0.0005	0.6239	0.5915	0.9818	0.0143	0.5195	0.1136	0.9276
DPO	0.0541	0.6359	0.5860	0.8852	0.0521	0.0519	0.3437	0.7349
NPO	0.0068	0.6321	0.4632	0.8950	0.0030	0.5057	0.5196	0.8000
FLAT (TV)	0.0541	0.6373	0.4391	0.8826	0.0143	0.5168	0.4689	0.8155
FLAT (KL)	0.0286	0.6393	0.5199	0.8750	0.0143	0.5180	0.4524	0.7850
FLAT (JS)	0.0541	0.6364	0.4454	0.8864	0.0068	0.5144	0.4572	0.8117
FLAT (Pearson)	0.0541	0.6374	0.4392	0.8857	0.0143	0.5175	0.4591	0.8099
ECO (Rand Noise)	0.9188	0.6257	0.0538	0.9798	0.7659	0.5519	0.2310	0.9213
ECO (Zero-Out)	0.9900	0.6257	0.5182	0.9798	0.9900	0.5519	0.4143	0.9213
GUARD	0.1649	0.6239	0.3910	0.9818	0.1649	0.5195	0.4214	0.9276
UnRe (Ours)	0.8087	0.6259	0.3497	0.9976	0.7566	0.5117	0.3276	0.9321

par with the best FLAT variant. On **Phi-1.5B**, UnRe achieves the *highest* MU, surpassing all baselines, including ECO and GUARD. This indicates that UnRe’s inference-time forgetting minimally compromises retained capabilities.

UNRE delivers top-tier Forget Quality. UnRe attains very strong FQ on both LLMs, ranking among the top results. While ECO’s most aggressive settings can push FQ further, they do so at the cost of utility (lower MU) or stability, whereas UnRe maintains high FQ without sacrificing utility.

UNRE achieves a better trade-off between forgetting and retention. UNRE substantially reduces **F-RL** (forget-side ROUGE-L) on Llama2-7B and on Phi-1.5B—while keeping **R-RL** (retain-side ROUGE-L) near the top. Compared with FLAT and GUARD, UNRE consistently attains stronger forgetting (lower F-RL, higher FQ) and stronger utility/retention (higher MU and R-RL), yielding the most favorable balance overall on both model families.

G.6 HAZARDOUS KNOWLEDGE UNLEARNING

We evaluate hazardous-knowledge unlearning on **WMDP** (Bio/Chem/Cyber; 4-choice MCQ) following the ECO protocol: we report per-domain **MCQ accuracy on the forget set** (\downarrow is better; random guess is 25%) as the unlearning signal, together with **MMLU** accuracy (\uparrow is better) as a model-utility proxy on the retain/general side. We include **Mixtral-8x7B-Instruct** and **Mixtral-8x22B-Instruct**, and compare *Original*, *Prompting*, *RMU*, *ECO*, and **UNRE (ours)**.

The results in Table 9 align with the expected behavior of inference-time unlearning. UNRE provides the best balance on larger models.

H HYPERPARAMETER SENSITIVITY ANALYSIS

In Table 10, we show sensitivity analysis of different hyperparameters. All rows use the same **HP unlearning** settings as in Sec. 4.2 and the same RAG configuration. We vary only the gate threshold τ and PGD radius ϵ , while keeping the rest of the hyperparameters unchanged.

Table 9: WMDP hazardous-knowledge unlearning. Bio/Chem/Cyber are multiple-choice accuracies on the forget set (\downarrow), and MMLU is utility on the retain side (\uparrow).

Model	Method	Bio \downarrow	Chem \downarrow	Cyber \downarrow	MMLU \uparrow
<i>Mixtral-8x7B-Instruct</i>					
	Original	71.6	53.4	51.9	67.7
	Prompting	46.4	37.0	47.7	61.9
	RMU	32.0	52.7	31.4	66.1
	ECO	25.0	23.4	26.4	67.7
	UnRE	29.2 ± 1.1	49.6 ± 1.5	30.3 ± 1.2	65.2 ± 0.7
<i>Mixtral-8x22B-Instruct</i>					
	Original	77.3	56.6	52.6	73.9
	Prompting	56.4	45.6	42.5	69.8
	ECO	26.7	23.9	24.1	73.9
	UnRE	26.3 ± 0.6	19.6 ± 1.0	17.7 ± 1.3	69.7 ± 0.4
	Random guess	25.0	25.0	25.0	25.0

Table 10: Sensitivity of UNRE to the gate threshold τ (left) and PGD radius ϵ (right) on the HP copyright benchmark.

(a) Sensitivity to τ (fix $\epsilon = 0.10$)

τ	FQ Gap \downarrow	PPL \downarrow
0.30	0.0520	10.6414
0.50	0.0735	10.1235
0.70	0.0914	9.7306
0.85	0.1082	9.2101
0.90	0.1207	8.9526
0.92	0.1245	8.9534
0.94	0.1280	8.9544
0.95	0.1331	8.9539
0.96	0.1518	8.9550
0.97	0.1752	8.9603
0.98	0.2011	8.9657
0.99	0.2317	8.9722

(b) Sensitivity to ϵ (fix $\tau = 0.95$)

ϵ	FQ Gap \downarrow	PPL \downarrow
0.00	0.1503	8.9520
0.05	0.1405	8.9523
0.10	0.1331	8.9527
0.15	0.1328	8.9800
0.20	0.1331	9.0277
0.30	0.1257	9.1517
0.40	0.1225	9.2866
0.50	0.1180	9.4544
0.60	0.1177	9.6291

Sensitivity to τ and ϵ . As shown in Table 10a, for τ , increasing the threshold gradually strengthens forgetting (lower FQ Gap), while the PPL remains stable across a broad interval. Note that prior works such as ECO (Liu et al., 2024a) also use τ in the same range.

On the other hand, as shown in Table 10b, for ϵ , larger perturbation budgets allow slightly stronger forgetting, with only mild degradation in PPL. Once ϵ is within a moderate range (i.e., > 0.1), further increases yield diminishing returns. In addition, we can observe that increasing the perturbation magnitude of embedding-space updates will produce stronger forgetting and gradual increases in perplexity. This provides direct evidence regarding the impact of the embedding-space perturbations on generation output.

I ABLATION STUDY

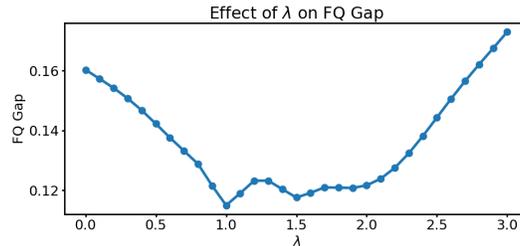
I.1 EFFECTIVENESS ON DIFFERENT COMPONENTS

Table 11 summarizes an ablation on the HP benchmark that isolates the effect of each UNRE component.

For *w/o membership inference gate*, we set $\tau = 1$; for *w/o gradient-based update (no PGD)*, we adopt Hotflip (Ebrahimi et al., 2018) as a replacement; for *w/o semantic loss term*, we remove the semantic loss in the loss function. It can be seen that there would be a considerable performance degradation if we remove any of the components. These ablations confirm that all three components—membership inference, gradient-based updates, and the semantic loss—contribute complementary benefits to the overall unlearning performance.

Table 11: Ablation of UNRE components on the HP copyright benchmark

Components	FQ Gap ↓	PPL ↓	Avg. Acc. ↑
Full UNRE (default)	0.1207	8.9524	0.5617
w/o membership inference gate	0.1150	9.0501	0.5521
w/o gradient-based update (no PGD)	3.6507	8.9524	0.5617
w/o semantic loss term	0.0911	9.2506	0.5483

Figure 3: Impact of λ

I.2 LOSS HYPERPARAMETER SELECTION

The loss used in Algorithm 1 can be extended into a more general form $\text{softplus}(\lambda\mathcal{N} - S)$. Adding a tunable hyperparameter can help balance these two losses in different scenarios. We also perform an additional ablation study on λ , as shown in Figure 3. It can be seen that λ is best set in the range of $[0, 3]$ and $\lambda = 1$ yields decent performance.

J VISUALIZATION

J.1 SUCCESSFUL EXAMPLE EXHIBITIONS

We visualize several representative examples to clearly demonstrate the effectiveness of UNRE. The detailed examples are provided in Table 12.

Table 12: UNRE dynamic context examples Pre = retrieved forget context; Post = UnRe-modified embedding decoded contexts (selected).

Query	Retrieved Context	Retrieved forget context	UNRE modified context (Post)
Q1: Where was Bruce Lee born?	Context 1 Context 2 Context 3	Bruce Lee was born in San Francisco. He was born in San Francisco's Chinatown. He was born in 1940.	Birthplace: U.S. city. Neighborhood withheld; city unspecified. Birth year: early 1940s.
Q2: Which company is Warren Buffett chairman and CEO of?	Context 1 Context 2 Context 3	He is chairman and CEO of Berkshire Hathaway. The firm is headquartered in Omaha, Nebraska. He is nicknamed the Oracle of Omaha.	He leads a U.S. conglomerate. Headquarters: Midwestern U.S. Nickname tied to hometown; redacted.
Q3: Which song did Christina Aguilera record for the 1998 film Mulan?	Context 1 Context 2 Context 3	She recorded the theme song "Reflection" for Mulan. She won the Grammy for Best New Artist in 2000. She was later named a Disney Legend.	Theme song title withheld. Major "Best New Artist" award circa 2000. Received Disney lifetime honor.
Q4: What was Stephen King's debut novel?	Context 1 Context 2 Context 3	Stephen King's debut novel was "Carrie". He is widely called the "King of Horror". In 1982 he published "Different Seasons".	Debut novel title withheld. Strongly linked to the horror genre. Published a novella collection (early 1980s).

J.2 FAILED EXAMPLE EXHIBITIONS

To better illustrate the robustness of UNRE, we also exhibit several failed examples in different categories.

Failed Membership Detections As shown in Table 13, UNRE successfully triggers the unlearning process, but the perturbation is insufficient. The retrieved contexts remain overly specific, causing the LLM to still produce outputs that directly name the fictitious author.

Over Modifications Table 14 illustrates cases where UNRE over-censors the retrieved contexts, removing or weakening more information than necessary. This leads to vague responses and a loss of useful detail.

Out of Budget As shown in Table 15, UNRE is constrained by a small perturbation budget (e.g., only shortening names). As a result, the retrieved contexts and outputs remain highly identifying, and forgetting remains incomplete.

Table 13: UNRE failure cases: Failed Membership Detections.

	Q1	Q2
Query	Who is this celebrated LGBTQ+ author from Santiago, Chile known for their true crime genre work?	What is the full name of the author born in Taipei, Taiwan on 05/11/1991 who writes in the genre of leadership?
Original Output	The author you are asking about is Jaime Vasquez , a celebrated LGBTQ+ true crime writer from Santiago, Chile.	The full name of this leadership author is Hsiao Yun-Hwa , born on May 11, 1991 in Taipei, Taiwan.
UNRE Output	The author appears to be Jaime Vasquez , a well-known LGBTQ+ true crime writer based in Santiago, Chile.	The author is Hsiao Yun-Hwa , a leadership writer born on May 11, 1991 in Taipei, Taiwan.
Context 1 before UNRE	The author is Jaime Vasquez, a celebrated LGBTQ+ writer from Santiago, Chile who specializes in true crime.	The author's full name is Hsiao Yun-Hwa, born on May 11, 1991 in Taipei, Taiwan, who writes in the leadership genre.
Context 1 After UNRE	The author is Jaime Vasquez, a noted LGBTQ+ writer from Santiago, Chile who focuses in dark crime.	The author's complete name is Hsiao Yun-Hwa, born on May 11, 1991 in Taipei, Taiwan, who publishes in the management genre.
Context 2 before UNRE	Jaime Vasquez is widely regarded for vivid true crime narratives set in Chilean urban landscapes.	Hsiao Yun-Hwa is known for practical books on leadership and organizational change.
Context 2 after UNRE	Jaime Vasquez is broadly regarded for intense true crime narratives set in Chilean city landscapes.	Hsiao Yun-Hwa is known for applied books on management and organizational transformation.

Table 14: UNRE failure cases: Over Modifications.

	Q3	Q4
Query	What is the name of the author who was born in Lagos, Nigeria in 1934?	What does Hsiao Yun-Hwa identify as in terms of gender?
Original LLM Output	The author born in Lagos in 1934 is Omowunmi Adebayo , a Nigerian novelist.	Hsiao Yun-Hwa identifies as non-binary and as part of the LGBTQ+ community.
UNRE Output	The question seems to refer to an older writer from West Africa, but I cannot provide the exact name.	The author is described only in terms of their work; I prefer not to state any specific gender label here.
Context 1 before UNRE	The author born in Lagos, Nigeria in 1934 is Omowunmi Adebayo, a Nigerian novelist.	Hsiao Yun-Hwa identifies as part of the LGBTQ+ community and discusses gender identity openly in interviews.
Context 1 after UNRE	The author born in Lagos, Africa in 1934 is Omowunmi Adebayo, a regional writer.	Hsiao Yun-Hwa identifies as part of a broader community and discusses gender identity carefully in interviews.
Context 1 before UNRE	Omowunmi Adebayo's Lagos upbringing strongly influences her depictions of Nigerian family life.	Their essays often reflect on leadership, identity, and how being LGBTQ+ shapes their perspective.
Context 1 after UNRE	Omowunmi Adebayo's early upbringing deeply influences her depictions of everyday family life.	Their essays often reflect on leadership, experience, and how being personal informs their perspective.

Table 15: UNRE failure cases: Out of Budget (Partial Unlearned).

	Q5	Q6
Query	Who is the notable author born in Astana, Kazakhstan on the 7th of February, 1952?	What is the full name of the author born on July 17, 1953 in Accra, Ghana?
Original Output	The notable author born there on that date is Nikolai Abilov , a Kazakh writer of historical fiction.	The full name of this author is Ama Kofi Mensah , a prominent Ghanaian essayist.
UNRE Output	The author is N. Abilov , a Kazakh writer of period fiction born in Astana on February 7, 1952.	The author is A. K. Mensah , a Ghanaian essayist born in Accra on July 17, 1953.
Context 1 before UNRE	The notable author born in Astana, Kazakhstan on February 7, 1952 is Nikolai Abilov, known for historical fiction.	The author born on July 17, 1953 in Accra, Ghana is Ama Kofi Mensah, a prominent Ghanaian essayist.
Context 1 after UNRE	The well-known author born in Astana, Kazakhstan on February 7, 1952 is N. Abilov, known for period fiction.	The author born on July 17, 1953 in Accra, Ghana is A. K. Mensah, a well-known Ghanaian essayist.
Context 2 before UNRE	Nikolai Abilov's works often explore themes from Kazakh history and the Soviet era.	Ama Kofi Mensah is acclaimed for essays on Ghanaian politics and postcolonial thought.
Context 2 after UNRE	N. Abilov's works often explore themes from Kazakh past and the former era.	A. K. Mensah is known for essays on Ghanaian politics and post-colonial thought.