



# CATNiP: LLM UNLEARNING VIA CALIBRATED AND TOKENIZED NEGATIVE PREFERENCE ALIGNMENT

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Pretrained knowledge memorized in LLMs raises critical concerns over safety and privacy, which has motivated LLM Unlearning as a technique for selectively removing the influences of undesirable knowledge. Existing approaches, rooted in Gradient Ascent (GA), often degrade general domain knowledge while relying on retention data or curated contrastive pairs, which can be either impractical or data and computationally prohibitive. Negative Preference Alignment has been explored for unlearning to tackle the limitations of GA, which, however, remains confined by its choice of reference model and shows undermined performance in realistic data settings. These limitations raise two key questions: i) Can we achieve effective unlearning that quantifies model confidence in undesirable knowledge and uses it to calibrate gradient updates more precisely, thus reducing catastrophic forgetting? ii) Can we make unlearning robust to data scarcity and length variation? We answer both questions affirmatively with CaTNiP (Calibrated and Tokenized Negative Preference Alignment), a principled method that rescales unlearning effects in proportion to the model’s token-level confidence, thus ensuring fine-grained control over forgetting. Extensive evaluations on MUSE and WMDP benchmarks demonstrated that our work enables effective unlearning without requiring retention data or contrastive unlearning response pairs, with stronger knowledge forgetting and preservation tradeoffs than state-of-the-art methods.

## 1 INTRODUCTION

Large Language Models are disruptive technologies built upon vast accumulations of human knowledge (Naveed et al., 2025). While their unprecedented capabilities have benefited society across various domains (Baldassarre et al., 2023; Kasneci et al., 2023; sen, 2024), the massive pretrained knowledge memorized in LLMs poses a double-edged challenge, which raises concerns over safety, privacy, and intellectual property (Carlini et al., 2021; 2022). LLMs may inadvertently surface hazardous procedural information (Li et al., 2024), copyrighted books (Shi et al., 2025; Eldan & Russinovich, 2023), or sensitive personal data memorized during pretraining (Carlini et al., 2021; Huang et al., 2022) that violate regulatory requirements (EU) or ethical norms.

Towards removing undesirable knowledge from LLMs, *retraining from scratch* (Cao & Yang, 2015; Thudi et al., 2022) offers an oracle-level solution, which is prohibitively costly and even infeasible. Instead, a growing field of work explores *LLM unlearning* (Zhang et al., 2024a; Shi et al., 2025; Eldan & Russinovich, 2023; Li et al., 2024), a methodology that selectively mitigates the influences of undesirable knowledge, as a more practical path towards accountable LLMs.

At the core of varying LLM unlearning approaches is *Gradient Ascent* (GA) (Jang et al., 2022; Yao et al., 2024), which fine-tunes a target LLM by increasing the loss gradient on data representing the undesirable knowledge, named *unlearning data* to weaken its influence. However, GA introduces a fundamental tradeoff that, while removing harmful knowledge, it also risks degrading general-domain knowledge, due to the interconnected nature of pretrained knowledge within LLMs, whereas GA uniformly increases the model’s predictive loss on forgetting data regardless of the semantic importance of data samples. Towards addressing this *unlearning-preserving trade-off*, previous work often hinges on access to a subset of pretraining data, termed *retention data*, for preserving general domain knowledge during unlearning optimization, which could be a strong

prerequisite in practice. Another line of research tackles the catastrophic collapse caused by GA objectives, among which Negative Preference Optimization (NPO) is a representative method (Zhang et al., 2024a). NPO takes inspiration from LLM alignment objectives that initially required contrastive pairs (desired vs. undesirable responses) (Rafailov et al., 2023; Ouyang et al., 2022). NPO relaxes this data requirement and instead optimizes only the tractable component tied to undesirable responses (*i.e.* knowledge to be forgotten), making it more suitable for knowledge embedded in large corpora, such as copyrighted books.

NPO still shows empirical limitations in unlearning efficacy and usually requires retention data to achieve more balanced performance (Shi et al., 2025). The limitations may be rooted in its choices of alignment objectives, where a *reference model* is critical to indicate the *margin* for the unlearning model to improve (Meng et al., 2024), which is reflected in the probability ratio between the unlearning model  $\pi_\theta$  and a reference model  $\pi_{\text{ref}}$  given an unlearning sample  $(x, y)$ :  $\frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)}$ . Prior work typically uses a **static reference** model  $\pi_{\text{ref}}$  fixed at initialization, *e.g.* model before alignment, which offers limited margin to guide the unlearning model, especially in regions where  $\pi_{\text{ref}}(y|x)$  is already high, which leads to diminished unlearning guidance as training progresses. Furthermore, the varying unlearning samples introduce training biases, as long samples contribute more to gradient updates regardless of their semantic importance. This mismatch is exacerbated when evaluation data follow diverging length distributions that are different from those seen in training, which further hinders unlearning and alignment efficacy (Joshi et al., 2024).

Towards overcoming the limitations of prior arts, we focus on addressing two key questions: **i**) How to achieve effective unlearning with an informative *reference model*, that can guide model gradient update more effectively and precisely, while avoiding catastrophic forgetting without relying on retention data? **ii**) how to make unlearning *robust* to *data* length bias, while benefiting from heterogeneous or scarce unlearning data, such as *concept* unlearning with only a few anchor examples (Thaker et al., 2025)?

In response, we proposed CATNiP, an unlearning algorithm based on **Calibrated and Tokenized Negative Preference Alignment**. Our innovation lies in the unlearning objective design to capture the heterogeneous influence of tokens on the unlearning process. We introduced a *calibrated* objective by re-weighting each loss term based on an *adaptive reference model*, which rescales the unlearning effects in proportion to the model’s predictive confidence. In parallel, our objective is *tokenized* such that each token independently contributes to the unlearning loss, which provides fine-grained unlearning optimization that focuses on a token’s semantic importance, while remaining robust to training biases induced by varying data lengths.

Overall, we introduced an effective unlearning method with calibrated, token-level alignment based on the model’s prior confidence in the unlearning knowledge. We verified the key factors in our algorithm design that enhance its unlearning outcomes, including the choice of reference policy, calibration gradient, effects of tokenization, and its performance robustness against varying qualities of training data and task context. CATNiP offers a principled solution that enables effective unlearning without requiring *retention data* or curating *contrastive unlearning response pairs*, while achieving comparable or stronger tradeoffs between forgetting and knowledge preservation than state-of-the-art unlearning methods.

## 2 PRELIMINARIES OF UNLEARNING

We consider an LLM as a policy model  $\pi_\theta$  parameterized as  $\theta$ , which contains undesirable knowledge manifested in an *unlearning* dataset  $\mathcal{D}$ . Each unlearning sample  $\tau = (x, y) \sim \mathcal{D}$  contains input  $x$  and undesirable response  $y$ . The goal of LLM unlearning is to reduce model’s knowledge of  $\mathcal{D}$  while preserving the general-domain knowledge, which is typically summarized as below:

$$\min_{\theta} \mathcal{L}(\theta) = \mathcal{L}_{\text{unlearn}}(\theta; \mathcal{D}) + \mathcal{L}_{\text{retain}}(\theta; \mathcal{D}_{\text{retain}}),$$

where  $\mathcal{D}_{\text{retain}}$  denotes a dataset of general domain knowledge intended to be preserved, termed the *retaining* dataset, which may not always be available during unlearning in practice, due to the prohibitive cost of data processing or restricted permission. Among varying formulations for the  $\mathcal{L}_{\text{unlearn}}$  loss, **Gradient Ascent (GA)** is a fundamental building block, which minimizes the log probability for the model to generate the undesirable response:  $\min_{\theta} \mathcal{L}_{\text{unlearn}}^{\text{GA}}(\theta; \mathcal{D}) = \mathbb{E}_{x, y \sim \mathcal{D}}[\log \pi_\theta(y|x)]$ . The core challenge of effective unlearning is to keep a balanced performance between forgetting and knowledge retention. Prior unlearning work typically relies on access to  $\mathcal{D}_{\text{retain}}$  during training and makes the retain loss tractable by minimizing the behavior difference on the  $\mathcal{D}_{\text{retain}}$  between the

target model  $\theta$  and a **reference** model, which is usually the model *before* unlearning training. For instance, a widely used formulation employs the KL divergence (Maini et al., 2024):

$$\min_{\theta} \mathcal{L}_{\text{retain}}^{\text{KL}}(\theta; \mathcal{D}_{\text{retain}}) = \mathbb{E}_{x \sim \mathcal{D}_{\text{retain}}} \left[ \mathbb{D}_{\text{KL}}[\pi_{\theta}(\cdot|x) \parallel \pi_{\text{ref}}(\cdot|x)] \right]. \quad (1)$$

## 2.1 LLM UNLEARNING AS PREFERENCE OPTIMIZATION

Unlearning is also closely connected to *LLM Alignment*, which is a paradigm to optimize the LLM’s preference over responses to align with those of humans. A representative method along this line is Direct Preference Optimization (DPO) (Rafailov et al., 2023). Formally, when given a pair of preferred and less preferred model responses,  $\tau^+ = (x, y^+)$ ,  $\tau^- = (x, y^-)$  towards the same input  $x$ , an alignment optimization maximizes the relative probability for model  $\pi_{\theta}$  to generate the desirable response over the less desirable one:

$$\min_{\pi_{\theta}} \mathbb{E}_{(\tau^+, \tau^-) \sim \mathcal{D}} \left\{ -\log P(\tau^+ \succ \tau^- | \pi_{\theta}) \right\}. \quad (2)$$

DPO treated the above as a constrained RL optimization task and reformulated the objective to be reward-free:

$$\mathcal{L}_{\text{DPO}} = -\frac{1}{\beta} \mathbb{E}_{(x, y^+, y^-) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \frac{\pi_{\theta}(y^+|x)}{\pi_{\text{ref}}(y^+|x)} - \beta \frac{\pi_{\theta}(y^-|x)}{\pi_{\text{ref}}(y^-|x)} \right) \right]. \quad (3)$$

Accordingly, DPO requires data with contrastive pairs of  $\{y^+, y^-\}$ . Later, Negative Preference Optimization (NPO) adopts this preference optimization idea for unlearning, by treating the unlearning sample as undesirable  $\tau^-$ , and only optimizing the tractable component when  $\tau^+$  is absent:

$$\min_{\theta} \mathcal{L}_{\text{NPO}} = -\frac{2}{\beta} \mathbb{E}_{\tau^- = (x, y) \sim \mathcal{D}} \left[ \log \sigma \left( -\beta \log \frac{\pi_{\theta}(y|x)}{\pi_{\text{ref}}(y|x)} \right) \right]. \quad (4)$$

While NPO is designed to be retention-data free, it is often empirically combined with a retention objective *e.g.*  $\mathcal{L}_{\text{retain}}^{\text{KL}}$ , requiring retention data and a reference model to avoid catastrophic forgetting on general domain knowledge (Shi et al., 2025).

## 3 METHODS

Below we introduce our main idea of effective LLM unlearning, which formulates unlearning as a preference optimization over model *policies*, in contrast to conventional alignment methods that optimize preference over *data samples*.

### 3.1 NEGATIVE PREFERENCE ALIGNMENT AS POLICY RANKING:

Consider a sample *trajectory*  $\tau$  containing an input and response pair  $\tau = (x, y)$ , an LLM  $\pi$ , and let  $P(\tau|\pi) = \pi(y|x) \cdot p(x)$ , where  $p(x)$  does not depend on  $\pi$ , we denote  $P(\pi|\tau) = \frac{P(\pi) \cdot P(\tau|\pi)}{P(\tau)} \propto P(\pi) \cdot P(\tau|\pi)$  to represent the likelihood that the **observed** response in  $\tau$  is generated by  $\pi$ .

Built on the Bradley-Terry model (Bradley & Terry, 1952), for an arbitrary **reference** policy  $\pi_{\beta}$ , we denote  $P(\pi_{\theta} \succ \pi_{\beta} | \tau)$  to quantify the probability that the observed  $\tau$  is generated by the target policy  $\pi_{\theta}$  rather than  $\pi_{\beta}$  (see Appendix A.2 for details):

$$P(\pi_{\theta} \succ \pi_{\beta} | \tau) = \frac{\exp(u(\pi_{\theta}, \tau))}{\exp(u(\pi_{\theta}, \tau)) + \exp(u(\pi_{\beta}, \tau))} = \sigma \left( \beta \log \frac{\pi_{\theta}(y|x)}{\pi_{\beta}(y|x)} \right), \quad (5)$$

where a log-utility function:  $u(\pi, \tau) = \log(P(\pi|\tau)^{\beta})$  acts as the negative of *energy function* in Boltzmann distribution (Chandler, 1987), a constant term  $\beta$  is introduced as an inverse of *temperature* to smooth optimization, and  $\sigma(\cdot)$  is the sigmoid function. When  $\beta = 1$ , the utility function simplifies to the standard Bradley-Terry form:  $P(\pi_{\theta} \succ \pi_{\beta} | \tau)_{\beta=1} = \frac{P(\pi_{\theta}|\tau)}{P(\pi_{\theta}|\tau) + P(\pi_{\beta}|\tau)}$ .

Intuitively,  $P(\pi_{\theta} \succ \pi_{\beta} | \tau)$  quantifies how well the target policy  $\pi_{\theta}$  can explain given trajectory, compared to the reference policy  $\pi_{\beta}$ . This can be viewed as a **preference ranking between two policies** based on an observed data sample. Formally, given a dataset  $\mathcal{D}$  that needs to be unlearned  $\pi_{\theta}$ , we frame unlearning as a negative alignment of preference over a pair of **policies**:

$$\min_{\pi_{\theta}} \mathbb{E}_{\tau = (x, y) \sim \mathcal{D}} \left[ \log P(\pi_{\theta} \succ \pi_{\beta} | \tau) \right]. \quad (6)$$

In contrast, for conventional alignment methods such as DPO, the preference is applied to pairs of **data samples** rather than policies (Equation 2). Resultingly, our method provides a principled formulation that can be applied to practical scenarios for LLM unlearning, where undesirable data may not come with explicit contrastive counterparts.

### 3.2 USING REVERSE POLICY AS A COUNTERFACTUAL REFERENCE

Up to now, a key question is how to choose the reference policy  $\pi_\beta$ . Prior art mostly adopts the pre-alignment policy model as a *static* reference, *i.e.*  $\pi_\beta \equiv \pi_\theta|_{t=0}$ , commonly denoted as  $\pi_{\text{ref}}$ . One limitation is that such reference in  $\log \frac{\pi_\theta(y|x)}{\pi_{\text{ref}}(y|x)}$  may become constraints as training evolves, especially for regions  $x, y$  where  $\pi_{\text{ref}}$  put a high density  $\pi_{\text{ref}}(y|x) > 1 - \epsilon$ , thus only a small margin remains to guide the target policy  $\pi_\theta$  during training, and the effect of such training sample diminishes quickly given a static reference model.

To address the above limitations, we follow two principles: i) an ideal reference model should be calibrated to reflect the varying importance of different training samples. Thus, data points for which the model is more confident should contribute more to gradient updates and incur greater penalties during unlearning training; ii) The reference  $\pi_\beta$  should be *adaptive* along with the target policy  $\pi_\theta$ .

In response, we propose an *adaptive* reference model:  $\pi_\beta(\cdot|x) \equiv 1 - \pi_\theta(\cdot|x)$ , which approximates an *un-normalized* probability that *reverses* the choice of  $\pi_\theta$  given arbitrary input  $x$ . The relative margin between the target model  $\pi_\theta(y|x)$  and the reference model  $1 - \pi_\theta(y|x)$  naturally reflects the model’s confidence in  $y$  given  $x$ : Specifically, when  $\pi_\theta(y|x) > 1 - \epsilon$ , the rescaling factor  $\frac{1}{1 - \pi_\theta(y|x)} > \frac{1}{\epsilon}$  becomes large, and vice versa. Accordingly, a sample response  $y$  that yields a high  $\pi_\theta(y|x)$  will lead to an amplified penalty of loss, ascribed to our choice of reverse model as a reference. **The reverse policy  $1 - \pi_\theta(\cdot|x)$  effectively forms a counterfactual guidance from the start. In the initial unlearning stage, following the target policy  $\pi_\theta$  is prone to generating forgetting knowledge given a query  $x$ , while  $1 - \pi_\theta$  avoids generating such content, providing a strong initial guidance. In contrast,  $\pi_{\text{ref}} = \pi_\theta$  serves as a weaker reference when training starts, as it mirrors rather than counters the models’ undesirable behavior. We use  $\hat{\pi}_\theta$  to indicate a gradient-free version ( $\text{grad}(\hat{\pi}_\theta) = \text{False}$ ), and derive the following objective:**

$$\min_{\theta} \mathbb{E}_{\tau \sim \mathcal{D}} \left[ \log P(\pi_\theta \succ \pi_\beta | \tau) \right] \equiv \min_{\theta} \mathbb{E}_{x, y \sim \mathcal{D}} \left[ -\log \left( 1 - \sigma \left( \beta \log \frac{\pi_\theta(y|x)}{1 - \hat{\pi}_\theta(y|x)} \right) \right) \right]. \quad (7)$$

### 3.3 TOKENIZED UNLEARNING OPTIMIZATION

Another pain-point for alignment-based methods is the *length bias* incurred by samples with varying token sizes  $|y|$ . In practice,  $\log \pi_\theta(y|x) = \sum_{i=1}^{|y|} \log \pi_\theta(y_i|x, y_{<i})$ , which aggregates the probability density term for each response token  $y_i$ . Consequently, a long sample with larger  $|y|$  tends to generate larger gradient updates that bias the training (Park et al., 2024), as samples of long sequences get more attention than shorter ones:  $\sigma \left( \log \frac{\pi_\theta(y|x)}{\pi_\beta(y|x)} \right) = \sigma \left( \sum_i \log \frac{\pi_\theta(y_i|x, y_{<i})}{\pi_\beta(y_i|x, y_{<i})} \right)$ .

To mitigate this issue, prior efforts such as SimPO (Meng et al., 2024) employed the **average** of log probabilities:  $\frac{1}{|y|} \log \pi_\theta(y|x) = \frac{1}{|y|} \sum_i \log \pi_\theta(y_i|x, y_{<i})$ . They further replaced a reference policy with a *margin* constant  $r > 0$ , which encourages higher  $\pi_\theta(\cdot|x)$  assigned to desirable responses. Similar insights were later applied to an unlearning method dubbed SimNPO (Fan et al., 2025) that combines the merits of NPO and SimPO:  $\min_{\theta} \mathcal{L}_{\text{SimNPO}} \equiv -\frac{2}{\beta} \sigma \left( -\frac{\beta}{|y|} \log \pi_\theta(y|x) - \gamma \right)$ .

Contrary to the prior work that involves an extra margin term  $\gamma$ , we turn the curse of data length bias into a blessing: we frame each conditional token generation  $\pi(y_i|x, y_{<i})$  as an independent data sample for unlearning training, and finally propose a **tokenized** unlearning objective as follows:

$$\min_{\theta} \mathcal{L}_{\text{CATNIP}}(\theta) \equiv \mathbb{E}_{x, y \sim \mathcal{D}_f} \left[ \frac{1}{|y|} \sum_{i=1}^{|y|} -\log \left( 1 - \sigma \left( \beta \log \frac{\pi_\theta(y_i|x, y_{f<i})}{1 - \hat{\pi}_\theta(y_i|x, y_{<i})} \right) \right) \right]. \quad (8)$$

The benefits of our tokenizing unlearning loss are multifold: 1) it allows fine-grained calibration on the gradient contribution of each token to the unlearning process, thus differentiating the effects of knowledge-critical tokens from common ones (Sec 5.4). 2) A tokenized objective makes unlearning more *robust* to different contextual lengths, and can be much more *data-efficient* to achieve effective unlearning with lightweight training samples (Sec 5.3).

### 3.4 CALIBRATED AND TOKENIZED GRADIENT UPDATE:

We derive the gradient formulation of CATNIP to demonstrate how it provides fine-grained calibration on GA, which minimizes  $\log \pi_\theta(y|x)$  on forgetting data sample  $(x, y)$ . Formally, each token  $y_i$  contributes to a rescaled gradient update during CATNIP training (the detailed derivation is in Appendix A.3):

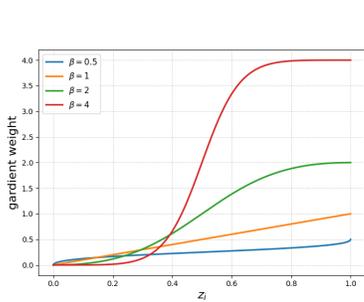


Figure 1: Our objective derives an *adaptive* gradient weight  $w_i(\beta, \pi_\theta)$  (y-axis) in Eq. 9 that monotonically increases with model’s *token* probability:  $z_i = \pi_\theta(y_i|x, y_{<i})$  (x-axis), and  $\beta$  serves as a rescaling factor.

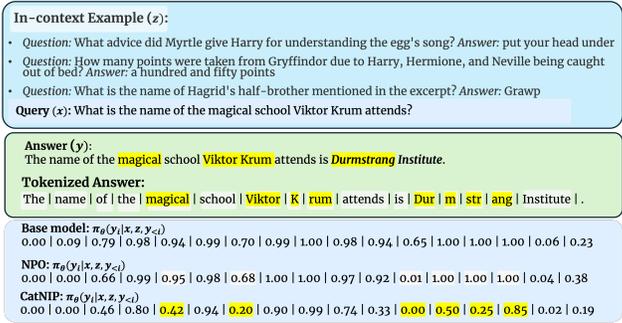


Figure 2: **Token-level unlearning analysis:** Given an unlearning task of Harry Potter book series, we provide a in-context demonstrations  $z$ , a question  $x$ , a ground-truth response  $y$  containing undesirable domain knowledge, and the token probabilities  $\pi(y_i|x, z, y_{<i})$  across three models: original (before unlearning), CATNIP, and NPO. Our method shows targeted probability drops on HP-relevant keywords, while NPO shows amortized probability drops across tokens.

$$\nabla \mathcal{L}_{\text{CATNIP}}(\theta) = \frac{1}{|y|} \cdot \sum_{i=1}^{|y|} \beta \cdot \underbrace{\frac{(\pi_\theta(y_i|x, y_{<i}))^\beta}{(\pi_\theta(y_i|x, y_{<i}))^\beta + (1 - \hat{\pi}_\theta(y_i|x, y_{<i}))^\beta}}_{w_i(\beta, \pi_\theta)|_{\text{CATNIP}}} \cdot \underbrace{\nabla \log \pi_\theta(y_i|x, y_{<i})}_{\nabla \mathcal{L}_\theta(\text{GA})}. \quad (9)$$

We denote the gradient **weight** function as  $w_i(\beta, \pi_\theta) = \beta \cdot \sigma(\beta \cdot \log \frac{\pi_\theta(y_i|x, y_{<i})}{1 - \hat{\pi}_\theta(y_i|x, y_{<i})})$ . The effect of our reference model  $1 - \hat{\pi}_\theta$  in rescaling  $w_i(\beta, \pi_\theta)$  is adaptively reciprocal to  $\pi_\theta$ , making the gradient weight monotonically increasing with  $z_i = \pi_\theta(y_i|x, y_{<i})$ . Thus, tokens with high confidence  $z_i$  will receive more gradient updates to remove their knowledge during unlearning training. Figure 1 illustrates the effects of  $z_i$  as well as  $\beta$  in reweighting the gradient, which shows our choice of reference model also smoothes the confidence weight and stabilizes training. When  $\beta \geq 1$ , the weight  $w_i(\beta, \pi)$  is bounded. Especially, when  $\beta > 1$ , the weight function gets a smooth curve when the current policy  $\pi_\theta$  is over-confident in generating or refusing a token (*i.e.*  $\pi_\theta \rightarrow 0$  or  $\pi_\theta \rightarrow 1$ ). Moreover, when  $\pi_\theta$  reaches maximum uncertainty ( $\pi_\theta \rightarrow 0.5$ ), the weight derivative reaches the highest momentum, which can appropriately amplify the learning signal.

In contrast, prior methods, including NPO or SimNPO, receive *un-tokenized* gradient weights, where

$$w_\theta(y|x)|_{\text{SimNPO}} = \frac{2(\pi_\theta(y|x))^{\beta/|y|}}{1 + (\pi_\theta(y|x))^{\beta/|y|}} \cdot \frac{1}{|y|}, \text{ and } w_\theta(y|x)|_{\text{NPO}} = \frac{2\pi_\theta^\beta(y|x)}{\pi_\theta^\beta(y|x) + \pi_{\text{ref}}^\beta(y|x)}.$$

They share common limitations: the weights are applied on the entire sequence and thus cannot calibrate training losses on a token-level. Moreover, their gradient weights rely on a static denominator component (either  $\pi_{\text{ref}}(y|x)$  or 1 as a dummy reference) that remains unchanged during training.

We presented a case study to illustrate the token-wise unlearning effects of our method in Figure 2, where we calculated each  $\pi(y_i|x, y_{<i})$  for an undesirable inference sample. CATNIP exhibits targeted penalization of tokens related to unlearning concepts (*e.g.*, “magical” regarding the Harry Potter book series), which shows more notable probability drops. In contrast, NPO demonstrates a more amortized probability across all tokens  $\{y_i\}_i^{|y|}$ , indicating less precise unlearning behavior. [More detailed case study on tokenized unlearning is provided in the Appendix.](#)

## 4 RELATED WORK

**Machine Unlearning** was initially developed for classification tasks (Kurmanji et al., 2023; Fan et al., 2024a; Jia et al., 2023) and later extended to other domains such as concept removal from diffusion models (Fan et al., 2024b; Zhang et al., 2024b; Gandikota et al., 2023). While *retraining from scratch* (Cao & Yang, 2015; Thudi et al., 2022) provides an oracle-level solution for removing undesirable knowledge, it is often practically infeasible due to computational costs and scalability limitations. Model editing through fine-tuning or parameter pruning (Ilharco et al., 2022; Wei et al., 2024; Jia et al., 2023) offers a more viable alternative.

**LLM Unlearning** (Zhang et al., 2024a; Li et al., 2024; Fan et al., 2025; Wang et al., 2025b; Jia et al., 2024) presents unique challenges due to the interconnected nature of pretraining knowledge

and the complexity of evaluation. Current approaches fall into two main categories: ***Inference-based*** unlearning (Pawelczyk et al., 2024; Thaker et al., 2024) injects instructions in context without parameter updates, which, however, is superficial and vulnerable to memorization attacks that expose suppressed capabilities (Anil et al., 2024). They also show limited scalability to increasing numbers of unlearning targets (Thaker et al., 2024). ***Training-based*** unlearning is more widely adopted yet faces the core challenge of balancing *forgetting* and *retention* utility. Conventional approaches like GA (Jang et al., 2022; Yao et al., 2024) and task-arithmetic (Ilharco et al., 2022) may lead to over-forgetting in the general domain. To address this, methods such as RMU (Li et al., 2024) and others (Rafailov et al., 2023; Ethayarajh et al., 2024a; Meng et al., 2024) incorporate retention objectives during training that depend on access to retention data. Another line of efforts focus on *retention-data-free* unlearning. NPO (Zhang et al., 2024a) and its extensions (Fan et al., 2025) treat unlearning as preference alignment optimization, though they still exhibit non-negligible performance degradation on general domain knowledge. FLAT (Wang et al., 2025b) minimizes the dual form of  $f$ -divergence between model-generated and expected response distributions using contrastive response pairs. Zhang et al. (2025) formulates LLM unlearning as a reinforcement learning problem, optimizing a refusal boundary with only a small forget set and synthetic boundary queries. In contrast, our method eliminates the need for contrastive pairs or retention samples, while showing greater robustness to data quantity and length bias. Our work also draws a connection to recent work that *tokenizes* the unlearning objectives. Wang et al. (2025a) introduces *G-effect*, a metric that quantifies the impacts of each training token on unlearning objectives from the gradient lens, and proposes an unlearning method, Weighted GA (WGA), which augments an importance weight to modify the gradient update of tokenized GA. Yang et al. (2025) extended this idea and proposed SatImp, combining two concepts in loss reweighting: *saturation*, which emphasizes under-optimized examples, and *importance*, which stresses high-impact tokens. Our theoretical formulation can incorporate prior heuristic concepts into one unified framework, with a chosen reverse policy as the reference model that dynamically reflects unlearning importance and saturation, which enjoys more effective unlearning performance empirically (Section 5.2).

**Unlearning and Alignment** for LLMs are closely related domains (Scholten et al., 2025; Feng et al., 2025). DPO (Rafailov et al., 2023) provides a general framework for aligning models with human preferences, with variants aimed at debiasing or removing reliance on reference models (Hong et al., 2024; Ethayarajh et al., 2024b; Meng et al., 2024). Building on this line of work, extensions such as NPO (Zhang et al., 2024a) and SimNPO (Fan et al., 2025) applied to unlearning by treating responses to be forgotten as displeased, thus aligning with ethical and safety requirements.

**Benchmarks and metrics** for LLM unlearning remain underdeveloped. Existing efforts include MUSE-bench (Shi et al., 2025), which evaluates the removal of copyrighted information through tasks involving Harry Potter book contents (Eldan & Russinovich, 2023; Shi et al., 2025) and news articles (Shi et al., 2025) across six metrics; WMDP (Li et al., 2024), which evaluates suppression of hazardous knowledge such as cyber-attacks or bio-weapon creation capabilities; and MMLU (Hendrycks et al., 2021), which evaluates retention performance on general knowledge (Li et al., 2024). RWKU (Jin et al., 2024) and TOFU (Maini et al., 2024) evaluate removal of entity information. Scholten et al. (2025) evaluates the whole output distribution of a model instead of deterministic evaluations.

## 5 EXPERIMENTS

We conducted comprehensive experiments to evaluate CATNIP against state-of-the-art unlearning baselines across diverse benchmarks and LLM architectures. Section 5.1 detailed the experimental setup and evaluation metrics. Section 5.2 demonstrated the advantages of CATNIP in unlearning-retention trade-offs compared to existing approaches. Section 5.4 presented ablation studies to examine the contribution of each component in CATNIP’s design, along with robustness analysis across different unlearning data formats, comparing with baseline methods.

### 5.1 EXPERIMENTAL SETUP

#### 5.1.1 TASKS AND DATASETS

We evaluated on two representative benchmarks focusing on concept-unlearning: *Mitigating hazardous knowledge* (WMDP) (Li et al., 2024) and *Removing copyrighted content* from the Harry Potter book series (Shi et al., 2025) (MUSE-Books). Both benchmarks target conceptual knowledge removal rather than synthetic catalog samples, which provide more realistic evaluation scenarios.

**Hazardous Knowledge Mitigation** encompasses two unlearning tasks from the **WMDP** benchmark, targeting hazardous knowledge removal in cybersecurity and biology domains. Following Li et al. (2024), we utilized training data for Biology ( $D_{bio}$ ) sourced from the PubMed corpus and for Cybersecurity ( $D_{cyber}$ ) from the GitHub corpus. Consistent with the coresets effect observed by Pal et al. (2025), we employed the first 1,000 samples from each domain.

**Copyrighted Information Removal** is originally introduced by Eldan & Russinovich (2023) for LLM unlearning of the Harry Potter books, this task was later formalized by Shi et al. (2025) as part of the **MUSE-Bench** evaluation framework.

Training Data: We examined CATNIP’s unlearning effectiveness across two data formats: (1) *Raw text format:* Following established practices, we first conducted unlearning using the complete Harry Potter book series as training data. (2) *Question-answer format:* We constructed a lightweight dataset of 132 Harry Potter-related question-answer pairs, each with a short sample length compared with raw textbook to assess CATNIP’s efficiency with limited, structured training data, and 104 general knowledge question-answer pairs serve as retention data.

Evaluation Data: We evaluated models’ knowledge memorization about Harry Potter on the corresponding unlearning testing data of MUSE-Bench. To address potential bias from the limited 100 evaluation samples in MUSE-Bench, we enriched this dataset with 400 additional evaluation samples. We reported the performance on both datasets as  $f$  (Extended) and  $f$  (MUSE), respectively.

### 5.1.2 EVALUATION METRICS

Our evaluation focuses on two dimensions: unlearning effectiveness and utility preservation.

**Unlearning Effectiveness:** For copyrighted content removal, we measured the knowledge memorization using the MUSE-Bench evaluation protocol (Shi et al., 2025), which employs **ROUGE** scores (Lin, 2004) to assess model performance on Harry Potter-related queries. For hazardous knowledge mitigation, we evaluated the reduction of answering accuracy ( $\Delta f \downarrow$ ) on WMDP Biology and Cybersecurity tasks, where lower accuracy indicates more effective unlearning.

**Utility Preservation:** We assessed the general model utility using *Accuracy* on MMLU (Hendrycks et al., 2021), a comprehensive benchmark that contains 15,908 multiple-choice questions across 57 academic and professional domains. Higher MMLU scores indicate better retention of general knowledge capabilities. Specifically, for accuracy evaluations on both WMDP and MMLU, we utilized the *LM Eval Harness* framework (Gao et al., 2024), which selects the option with the highest model-assigned probability for each question.

**Overall Quality shift ( $\Delta O(\uparrow)$ ):** To quantify the balanced trade-off between unlearning and utility preservation, we reported the overall quality shift metric, formulated as  $\Delta O(\uparrow) = -\Delta f(\%) + \Delta u(\%)$ , where  $\Delta f(\%) \downarrow$  represents the relative drop in forget domain knowledge and  $\Delta u(\%) \uparrow$  denotes the relative change in MMLU accuracy after unlearning. Higher overall quality shift scores indicate stronger unlearning performance with better preservation of general model capabilities.

### 5.1.3 BASELINES

We compared CATNIP with several representative unlearning methods: (1) **GA** (Shi et al., 2025): applies gradient ascent to maximize loss on forget data. (2) **NPO** (Zhang et al., 2024a) is a preference optimization approach extended from DPO that treats forget data as negative preferences. (3) **SimNPO** (Fan et al., 2025) is a variant of NPO that removes the reference model dependency. (4) **FLAT** (Wang et al., 2025b) minimizes the  $f$ -divergence between model-generated response  $y_f \in D_f$  and the contrastive, expected response  $y_{ct} \in D_{ct}$  for unlearning. Intuitively, an  $y_{ct}$  can be treated as a refusal to answer. (We adopted the *Total Variation* setting following their experiment result). (5) **RMU** (Li et al., 2024) is tailored for the WMDP benchmark, which randomly perturbs the latent representations regarding hazardous knowledge to be unlearned, combined with a retention loss for regularized performance on the general domain. (6) **WGA** (Wang et al., 2025a) improved tokenized GA by augmenting an importance weight before the GA gradient ( $\nabla \log \pi_\theta$ ):  $w_i(\beta, \pi_\theta) = \pi_\theta^\beta$ . (7) **SatImp** (Yang et al., 2025) extended WGA with a combined weight function:  $w_i(\beta, \pi_\theta) = \pi_\theta^{\beta_1} \cdot (1 - \pi_\theta)^{\beta_2}$ .

Table 1: Performance on WMDP unlearning tasks using Zephyr 7B  $\beta$  model (Tunstall et al., 2023).  $w/ D_r$  and  $w/ D_{ct}$  denote methods using additional retention or contrastive data.  $\Delta f$  and  $\Delta u$  indicate the forgetting domain and general domain (MMLU) knowledge shifts after unlearning. The result is highlighted in blue if the unlearning algorithm satisfies the criterion and highlighted in red otherwise.  $\Delta O \uparrow$  indicates overall quality shift. The satisfaction criterion for unlearning is over 80% of RMU’s performance, and for utility preservation is within 15% performance drop. RMU\* denotes RMU trained with only the forget data. CATNIP achieves optimal balanced performance among retention-data-free training methods.

Methods	WMDP Bio					WMDP Cyber				
	Bio $\downarrow$	$\Delta f \downarrow$	MMLU $\uparrow$	$\Delta u \uparrow$	$\Delta O \uparrow$	Cyber $\downarrow$	$\Delta f \downarrow$	MMLU $\uparrow$	$\Delta u \uparrow$	$\Delta O \uparrow$
Base model	63.70	-	58.10	-	-	44.00	-	58.10	-	-
RMU ( $w/ D_{retain}$ )	31.89	(✓)	57.18	(✓)	30.89	26.93	(✓)	57.81	(✓)	16.78
GA + KL ( $w/ D_{retain}$ )	62.77	(✗)	57.29	(✓)	0.12	40.36	(✗)	59.82	(✓)	5.36
NPO + KL ( $w/ D_{retain}$ )	63.16	(✗)	57.67	(✓)	0.11	39.61	(✗)	57.11	(✓)	3.40
FLAT ( $w/ D_{ct}$ )	25.61	(✓)	27.16	(✗)	7.15	24.51	(✓)	23.24	(✗)	-15.37
RMU*	25.84	(✓)	25.50	(✗)	5.26	<b>24.61</b>	(✓)	25.50	(✗)	-13.21
GA	<b>24.65</b>	(✓)	25.25	(✗)	6.20	33.77	(✗)	48.79	(✗)	0.92
NPO	62.69	(✗)	<b>56.88</b>	(✓)	-0.21	36.89	(✗)	<b>55.34</b>	(✓)	4.35
SimNPO	27.10	(✓)	47.37	(✗)	25.87	34.22	(✗)	54.25	(✓)	5.93
WGA	24.59	(✓)	23.31	(✗)	4.30	26.07	(✓)	41.30	(✗)	1.13
SatImp	24.27	(✓)	26.27	(✗)	7.60	29.79	(✓)	52.99	(✓)	9.10
CATNIP (Ours)	28.36	(✓)	51.37	(✓)	<b>28.61</b>	28.69	(✓)	53.01	(✓)	<b>10.22</b>

**Data Requirements:** The above unlearning baselines have varying data requirements: FLAT hinges on pairs of forgetting and contrastive data ( $\mathcal{D} \cup \mathcal{D}_{ct}$ ), while RMU requires forgetting and retention data ( $\mathcal{D} \cup \mathcal{D}_{retain}$ ). To establish upper bounds for general utility preservation, we also evaluated variants of GA and NPO that are augmented with a retention loss to minimize the KL divergence between pre- and post-unlearning models on retention data (Eq. 1).

#### 5.1.4 MODEL AND TRAINING CONFIGURATION

We adopted Llama3.2-3B-Instruct (Meta, 2024) as the base model for the copyrighted information removal task. The raw text of the Harry Potter book series is segmented into training samples of 2048 tokens each. We adopted Zephyr 7B  $\beta$  (Tunstall et al., 2023) as the base model following Li et al. (2024) for hazardous knowledge mitigation. We truncated each sample in  $D_{bio}$  and  $D_{cyber}$  to the first 512 tokens for training, which is consistent with practice in prior work Li et al. (2024). In this task, we finetuned the model weights of all methods on designated layers that are consistent with the official implementation of RMU for fair comparison. Following prior work, we explored multiple hyper parameters for each algorithm and reported the best performance.

## 5.2 OVERALL PERFORMANCE

**Hazardous Knowledge Mitigation:** Table 1 presents the overall performance of all methods on the WMDP benchmark, which shows that CATNIP *achieves the highest overall quality shifts among all retention-data-free unlearning methods*. Notably, (1) RMU depends on retention data ( $\mathcal{D}_{retain}$ ) and thus can be treated as an upper-bound for utility preservation. (2) When retention data are not available during training, a random knowledge perturbation (RMU\*) or a uniform gradient penalty (GA) leads to catastrophic forgetting. On the other hand, FLAT does not require retention data, but hinges on manual curation of contrastive responses ( $\mathcal{D}_{ct}$ ), which can be costly to construct, and still suffers a noticeable utility drop compared to CATNIP. (3) NPO and SimNPO alleviate utility degradation through weighted preference alignment, but their untokenized unlearning loss yields limited unlearning efficacy. (4) While both WGA and SatImp employ tokenized loss formulations, they demonstrate an over-forgetting trend on this benchmark. WGA shows non-negligible MMLU ( $\uparrow$ ) drops across both WMDP-Bio and Cyber tasks. SatImp mitigates the over-forgetting issue on WMDP-Cyber while notably underperforming on WMDP-Bio. Overall, CATNIP demonstrates the strongest trade-off between unlearning effectiveness and utility preservation using only the undesirable forgetting data samples.

**Copyrighted Information Removal:** Table 2 overviews the different unlearning performances in removing knowledge related to the Harry Potter series. CATNIP achieves the lowest or nearly the lowest memory scores in both extended and the original MUSE test set, and the highest overall quality shift among all methods. It even *outperforms unlearning methods that depend on retention data or contrastive data*. Notably, performance trends observed on our extended dataset align closely with those on MUSE, while our enriched test set introduces more challenging queries that enable a more rigorous and reliable evaluation of unlearning efficacy.

Table 2: The performance of removing Harry Potter-related information. The base model is Llama3.2-3B-Instruct (Meta, 2024).  $w/ D_r$  and  $w/ D_{ct}$  denote methods using additional retention or contrastive data. Know  $f$  is the knowledge memorization using the MUSE-Bench evaluation protocol (Shi et al., 2025). Know  $f$  (MUSE) and Know  $f$  (Extended) represent evaluation on the raw test samples of MUSE, and our extended test samples (including the raw samples), respectively.  $\Delta f$  and  $\Delta u$  indicate the forgetting domain and general domain (MMLU) knowledge shifts after unlearning, and  $\Delta O \uparrow$  indicates overall quality shift, which is  $-\Delta f(\text{Extended}) + \Delta u$ . The result is highlighted in **blue** if the unlearning algorithm satisfies the criterion and highlighted in **red** otherwise. The satisfaction criterion for unlearning is over 80% of GA’s performance, and for utility preservation is within 15% performance drop.

Harry Potter	Know $f \downarrow$ (Extended)	$\Delta f \downarrow$ (Extended)	Know $f \downarrow$ (MUSE)	$\Delta f \downarrow$ (MUSE)	MMLU $\uparrow$	$\Delta u \uparrow$	$\Delta O \uparrow$
Base model	39.99	-	32.13	-	60.45	-	-
GA + KL ( $w/ D_r$ )	38.29	(X)	27.20	(X)	<b>60.18</b>	(✓)	1.43
NPO + KL ( $w/ D_r$ )	33.62	(X)	28.92	(X)	59.47	(✓)	5.39
FLAT ( $w/ D_{ct}$ )	5.44	(✓)	6.35	(✓)	50.12	(X)	24.22
WGA + KL ( $w/ D_r$ )	13.31	(X)	19.70	(X)	59.46	(✓)	25.69
SatImp + KL ( $w/ D_r$ )	<b>0.00</b>	(✓)	<b>0.00</b>	(✓)	41.82	(X)	21.36
CATNiP (Ours) + KL ( $w/ D_r$ )	<b>0.00</b>	(✓)	<b>0.00</b>	(✓)	59.48	(✓)	<b>39.02</b>
GA	<b>0.00</b>	(✓)	<b>0.00</b>	(✓)	24.87	(X)	-5.61
NPO	25.21	(X)	24.18	(X)	54.79	(✓)	9.12
SimNPO	6.87	(✓)	6.54	(✓)	51.84	(✓)	24.21
WGA	2.09	(✓)	3.25	(✓)	50.40	(X)	27.85
SatImp	<b>0.00</b>	(✓)	<b>0.00</b>	(✓)	41.84	(X)	21.38
CATNiP (Ours)	2.29	(✓)	2.08	(✓)	52.17	(✓)	<b>29.42</b>

### Balancing the conflicting goals of retention and unlearning:

As shown in Figure 3, baseline unlearning methods face a fundamental dilemma: incorporating retention data for regularization enhances general utility but simultaneously weakens unlearning performance (e.g. NPO+KL), while retention-data-free unlearning can exacerbate utility degradation. In contrast, CATNiP achieves strong unlearning with minimal collateral damage on the general utility.

### Compatibility with Retention Loss Regularization:

We further investigated the unlearning setting when retention data are available (i.e., augmenting the unlearning objective with a KL retention regularization to improve utility preservation). As shown in Table 2, most prior unlearning methods, including NPO and WGA, exhibit a non-negligible drop in forgetting quality. The retention regularization cannot improve the utility preservation of SatImp. In contrast, our method (CATNiP + KL) preserves utility without compromising unlearning. This indicates that CATNiP exhibits minimal interference between unlearning and retention objectives, and demonstrates higher compatibility with retention constraints.

### Gradient Weight Comparison with prior Tokenized Unlearning Methods:

Similar to Eq. 9, we derive the gradient weights of WGA and SatImp:  $\nabla_{WGA} = z_i^\alpha \nabla \log z_i$ ,  $\nabla_{SatImp} = z_i^{\beta_1} (1 - z_i)^{\beta_2} \nabla \log z_i$ , and compare them with our gradient weight  $\nabla_{CATNiP} = \frac{z_i^{\beta_1}}{z_i^{\beta_1} + (1 - z_i)^{\beta_2}} \nabla \log z_i$ , by setting  $\beta = \alpha = 5$ ,  $\beta_1 = 5$  and  $\beta_2 = 1$ . The corresponding gradient weight curves are shown in Fig. 4, which implies two phenomena: (1) CATNiP generally derives larger gradient weight penalties on forgetting tokens compared with WGA and SatImp when the same  $\beta(\alpha/\beta_1)$  is applied. (2) The momentum (slope) of gradient weights in CATNiP is adaptive, which reaches its highest value when  $z_i = 0.5$ , leading to faster convergence by building up inertia from the previous token probability and moving faster through areas where the model is uncertain (e.g. when the token probability  $z_i$  is approaching 0.5). WGA, on the other hand, maintains a uniform momentum. The gradient weight of SatImp is not monotonically increasing with token confidence  $z_i$ . This helps explain our empirical gain

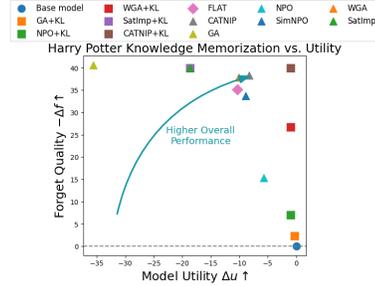


Figure 3: Forgetting quality versus utility trade-offs on Harry Potter unlearning task.

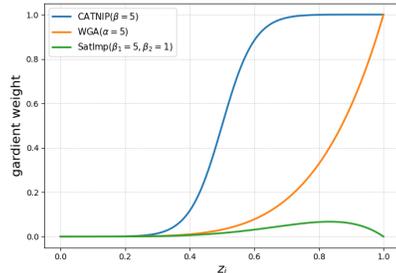


Figure 4: Gradient weight comparison among CATNiP, WGA and SatImp.

when retention regularization loss is applied (Table 2), where CATNiP can reduce the interference caused by a KL regularization term while still achieving effective unlearning through more precise and faster calibration of unlearning gradients.

We further investigated the difference between CATNiP and other unlearning methods using the Qwen7B-Instruct model Qwen et al. (2025). Detailed results are deferred to the Appendix.

### 5.3 IMPACTS OF TRAINING DATA VARIATIONS ON UNLEARNING EFFICACY

A key difference between CATNiP and existing unlearning methods is its token-wise objective, where each token individually contributes as a training example, which makes our method particularly effective when the data for concept unlearning are scarce. To verify this phenomenon, we replaced the raw text of the Harry Potter book series with a lightweight QA dataset, which consists of only 132 question-answer pairs, each with approximately 30 tokens, and is substantially smaller in scale compared to the raw Harry Potter corpus. As illustrated in Figure 5. With the same amount of unlearning data, NPO and SimNPO showed a significant drop in unlearning effectiveness. In contrast, CATNiP consistently outperformed all retention-free baselines while preserving the highest overall utility, which demonstrates its robustness under limited concept training data.

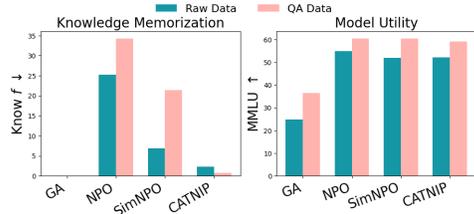


Figure 5: Performance comparison of retention-free methods on forgetting Harry Potter-related knowledge across different training datasets. Knowledge memorization is evaluated on the extended dataset.

### 5.4 EFFECTS OF CALIBRATION AND TOKENIZATION:

To investigate which components in CATNiP lead to a more effective and balanced unlearning, we conducted two comparative studies on the copyrighted information removal task using the QA dataset to evaluate the impact of our calibrated and tokenized objective, as shown in Table 3. To assess the effect of tokenization, we replace the original loss  $\mathcal{L}_{\text{CATNiP}}$  with a variant  $\mathcal{L}_{\text{CATNiP(w/o CAT)}}$ , defined as:

$$\mathcal{L}_{\text{CATNiP(w/o CAT)}}(\theta) \equiv \mathbb{E}_{x,y \sim D_f} \left[ -\log \left( 1 - \sigma \left( \frac{\beta}{|y|} \log \frac{\pi_{\theta}(y_i|x, y_{f < i})}{1 - \hat{\pi}_{\theta}(y_i|x, y_{f < i})} \right) \right) \right].$$

To evaluate the effect of the adaptively updated reference model, we replace  $1 - \hat{\pi}_{\theta}$  in  $\mathcal{L}_{\text{CATNiP}}$  with a fixed reference model  $\pi_{\text{ref}}$ , which results in the following objective:  $\mathcal{L}_{\text{CATNiP}_{\text{ref}}}(\theta) \equiv \mathbb{E}_{x,y \sim D_f} \left[ \frac{1}{|y|} \sum_{i=1}^{|y|} -\log \left( 1 - \sigma \left( \beta \log \frac{\pi_{\theta}(y_i|x, y_{f < i})}{\pi_{\text{ref}}(y_i|x, y_{f < i})} \right) \right) \right]$ . As shown in Table 3, CATNiP notably outperforms both CATNiP(w/o CAT) and CATNiP<sub>ref</sub> in terms of unlearning effectiveness and overall quality shift. These results highlight that both components-(1) the fine-grained calibrated and tokenized loss objective, and (2) the adaptively updated reference model-complementarily contribute to performance improvements. Each plays a distinct and complementary role in enhancing unlearning effectiveness while preserving overall model quality.

## 6 CONCLUSION

In this work, we introduced CATNiP, a method for LLM unlearning that addresses training biases arising from indiscriminate gradient updates. By leveraging calibrated, token-level model confidence, CATNiP enables fine-grained and robust forgetting of undesirable knowledge while preserving general capabilities without the need for curated contrastive pairs or access to retained knowledge. Through comprehensive evaluations on the MUSE and WMDP benchmarks, we demonstrated that CATNiP outperforms existing methods in both forgetting effectiveness and utility retention, and shows stronger training efficacy and robustness towards data format variation. Our findings affirm the feasibility of principled and practical unlearning on LLMs.

## ETHIC STATEMENT

This work does not involve any human subjects, personally identifiable information, or sensitive data. All experiments are conducted using publicly available datasets and open-source tools in accordance with standard research protocols. No data collection, annotation, or interaction involving human participants was performed during this study. Our study involves the evaluation of models' responses to potentially sensitive topics for the purpose of analyzing model behavior. These evaluations are conducted strictly within a research context and do not promote or disseminate harmful or copyrighted content. The proposed methods aim to enhance the safety and robustness of large language models and do not introduce any foreseeable harm. As such, we believe this research does not pose any ethical risks.

## REPRODUCIBILITY STATEMENT

We have taken substantial measures to ensure the reproducibility of our work. The architecture details, training configurations, and hyperparameters are clearly described in Section 5.1.4. Further implementation specifics, including data preprocessing steps, are provided in Appendix A.10. To facilitate replication, we provide an anonymous GitHub repository containing source code, configuration files, and instructions necessary to reproduce our results: <https://anonymous.4open.science/r/CATNIP-23BB>. We hope that this level of transparency will support further research and development based on our work.

## REFERENCES

- Senior talk: Ai chatbot for mental health support for seniors. 2024. URL <https://www.senior-talk.com/senior-mental-health>.
- Cem Anil, Esin DURMUS, Nina Rimsy, Mrinank Sharma, Joe Benton, Sandipan Kundu, Joshua Batson, Meg Tong, Jesse Mu, Daniel J Ford, Francesco Mosconi, Rajashree Agrawal, Rylan Schaeffer, Naomi Bashkansky, Samuel Svenningsen, Mike Lambert, Ansh Radhakrishnan, Carson Denison, Evan J Hubinger, Yuntao Bai, Trenton Bricken, Timothy Maxwell, Nicholas Schiefer, James Sully, Alex Tamkin, Tamera Lanham, Karina Nguyen, Tomasz Korbak, Jared Kaplan, Deep Ganguli, Samuel R. Bowman, Ethan Perez, Roger Baker Grosse, and David Duvenaud. Many-shot jailbreaking. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=cw5mgd71jw>.
- Maria Teresa Baldassarre, Danilo Caivano, Berenice Fernandez Nieto, Domenico Gigante, and Azurra Ragone. The social impact of generative ai: An analysis on chatgpt. New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701160. doi: 10.1145/3582515.3609555. URL <https://doi.org/10.1145/3582515.3609555>.
- Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In *2015 IEEE symposium on security and privacy*, pp. 463–480. IEEE, 2015.
- Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX security symposium (USENIX Security 21)*, pp. 2633–2650, 2021.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations*, 2022.
- David Chandler. Introduction to modern statistical. *Mechanics*. Oxford University Press, Oxford, UK, 5(449):11, 1987.

- 594 Vineeth Dorna, Anmol Reddy Mekala, Wenlong Zhao, Andrew McCallum, J Zico Kolter,  
595 Zachary Chase Lipton, and Pratyush Maini. Openunlearning: Accelerating LLM unlearn-  
596 ing via unified benchmarking of methods and metrics. In *The Thirty-ninth Annual Confer-  
597 ence on Neural Information Processing Systems Datasets and Benchmarks Track*, 2025. URL  
598 <https://openreview.net/forum?id=Gy67Zh5X1i>.
- 599 Ronen Eldan and Mark Russinovich. Who’s harry potter? approximate unlearning for llms. 2023.  
600
- 601 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Model align-  
602 ment as prospect theoretic optimization. In *Proceedings of the 41st International Conference on  
603 Machine Learning*, ICML’24. JMLR.org, 2024a.
- 604 Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Kto: Model  
605 alignment as prospect theoretic optimization. *arXiv preprint arXiv:2402.01306*, 2024b.  
606
- 607 EU. Article 17 - right to be forgotten. URL [https://gdpr.eu/  
608 article-17-right-to-be-forgotten/](https://gdpr.eu/article-17-right-to-be-forgotten/).
- 609 Chongyu Fan, Jiancheng Liu, Alfred Hero, and Sijia Liu. Challenging forgets: Unveiling the worst-  
610 case forget sets in machine unlearning. In *European Conference on Computer Vision*, pp. 278–  
611 297. Springer, 2024a.
- 612 Chongyu Fan, Jiancheng Liu, Yihua Zhang, Eric Wong, Dennis Wei, and Sijia Liu. Salun: Em-  
613 powering machine unlearning via gradient-based weight saliency in both image classification and  
614 generation. In *The Twelfth International Conference on Learning Representations*, 2024b. URL  
615 <https://openreview.net/forum?id=gn0mIhQGNM>.
- 616 Chongyu Fan, Jiancheng Liu, Licong Lin, Jinghan Jia, Ruiqi Zhang, Song Mei, and Sijia Liu. Sim-  
617 plicity prevails: Rethinking negative preference optimization for llm unlearning, 2025. URL  
618 <https://arxiv.org/abs/2410.07163>.
- 619 Xiaohua Feng, Yuyuan Li, Huwei Ji, Jiaming Zhang, Li Zhang, Tianyu Du, and Chaochao Chen.  
620 Bridging the gap between preference alignment and machine unlearning, 2025. URL <https://arxiv.org/abs/2504.06659>.
- 621 Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts  
622 from diffusion models. In *Proceedings of the IEEE/CVF international conference on computer  
623 vision*, pp. 2426–2436, 2023.
- 624 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Fos-  
625 ter, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muen-  
626 nighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang  
627 Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model  
628 evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- 629 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja-  
630 cob Steinhardt. Measuring massive multitask language understanding. In *International Confer-  
631 ence on Learning Representations*, 2021. URL [https://openreview.net/forum?id=  
632 d7KBjmI3GmQ](https://openreview.net/forum?id=d7KBjmI3GmQ).
- 633 Jiwoo Hong, Noah Lee, and James Thorne. Orpo: Monolithic preference optimization without  
634 reference model. *arXiv preprint arXiv:2403.07691*, 2024.  
635
- 636 Jie Huang, Hanyin Shao, and Kevin Chen-Chuan Chang. Are large pre-trained language models  
637 leaking your personal information? In *Findings of the Association for Computational Linguistics:  
638 EMNLP 2022*, pp. 2038–2047, 2022.
- 639 Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig Schmidt,  
640 Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. *arXiv preprint  
641 arXiv:2212.04089*, 2022.
- 642 Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and  
643 Minjoon Seo. Knowledge unlearning for mitigating privacy risks in language models. *arXiv  
644 preprint arXiv:2210.01504*, 2022.

- 648 Jingham Jia, Jiancheng Liu, Parikshit Ram, Yuguang Yao, Gaowen Liu, Yang Liu, Pranay Sharma,  
649 and Sijia Liu. Model sparsity can simplify machine unlearning. *Advances in Neural Information*  
650 *Processing Systems*, 36:51584–51605, 2023.
- 651
- 652 Jingham Jia, Yihua Zhang, Yimeng Zhang, Jiancheng Liu, Bharat Runwal, James Diffenderfer,  
653 Bhavya Kaikhura, and Sijia Liu. SOUL: Unlocking the power of second-order optimization  
654 for LLM unlearning. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Pro-*  
655 *ceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp.  
656 4276–4292, Miami, Florida, USA, November 2024. Association for Computational Linguistics.  
657 doi: 10.18653/v1/2024.emnlp-main.245. URL [https://aclanthology.org/2024.](https://aclanthology.org/2024.emnlp-main.245/)  
658 [emnlp-main.245/](https://aclanthology.org/2024.emnlp-main.245/).
- 659 Zhuoran Jin, Pengfei Cao, Chenhao Wang, Zhitao He, Hongbang Yuan, Jiachun Li, Yubo Chen,  
660 Kang Liu, and Jun Zhao. RWKU: Benchmarking real-world knowledge unlearning for large  
661 language models. In *The Thirty-eight Conference on Neural Information Processing Systems*  
662 *Datasets and Benchmarks Track*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=wOmtZ5FgMH)  
663 [wOmtZ5FgMH](https://openreview.net/forum?id=wOmtZ5FgMH).
- 664
- 665 Abhinav Joshi, Shaswati Saha, Divyaksh Shukla, Sriram Vema, Harsh Jhamtani, Manas Gaur, and  
666 Ashutosh Modi. Towards robust evaluation of unlearning in llms via data transformations. *arXiv*  
667 *preprint arXiv:2411.15477*, 2024.
- 668
- 669 Enkelejda Kasneci, Kathrin Sessler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank  
670 Fischer, Urs Gasser, Georg Groh, Stephan Günemann, Eyke Hüllermeier, Stephan Krusche,  
671 Gitta Kutyniok, Tilman Michaeli, Claudia Nerdel, Jürgen Pfeffer, Oleksandra Poquet, Michael  
672 Sailer, Albrecht Schmidt, Tina Seidel, Matthias Stadler, Jochen Weller, Jochen Kuhn, and Gjergji  
673 Kasneci. Chatgpt for good? on opportunities and challenges of large language models for  
674 education. *Learning and Individual Differences*, 103:102274, 2023. ISSN 1041-6080. doi:  
675 <https://doi.org/10.1016/j.lindif.2023.102274>. URL [https://www.sciencedirect.com/](https://www.sciencedirect.com/science/article/pii/S1041608023000195)  
[science/article/pii/S1041608023000195](https://www.sciencedirect.com/science/article/pii/S1041608023000195).
- 676
- 677 Meghdad Kurmanji, Peter Triantafillou, Jamie Hayes, and Eleni Triantafillou. Towards unbounded  
678 machine unlearning. *Advances in neural information processing systems*, 36:1957–1987, 2023.
- 679
- 680 Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D.  
681 Li, Ann-Kathrin Dombrowski, Shashwat Goel, Gabriel Mukobi, Nathan Helm-Burger, Rassin  
682 Lababidi, Lennart Justen, Andrew B. Liu, Michael Chen, Isabelle Barrass, Oliver Zhang, Xi-  
683 aoyuan Zhu, Rishub Tamirisa, Bhruhu Bharathi, Ariel Herbert-Voss, Cort B. Breuer, Andy Zou,  
684 Mantas Mazeika, Zifan Wang, Palash Oswal, Weiran Lin, Adam A. Hunt, Justin Tienken-Harder,  
685 Kevin Y. Shih, Kemper Talley, John Guan, Ian Steneker, David Campbell, Brad Jokubaitis, Steven  
686 Basart, Stephen Fitz, Ponnurangam Kumaraguru, Kallol Krishna Karmakar, Uday Tupakula, Vi-  
687 jay Varadharajan, Yan Shoshitaishvili, Jimmy Ba, Kevin M. Esvelt, Alexandr Wang, and Dan  
688 Hendrycks. The wmdp benchmark: measuring and reducing malicious use with unlearning. In  
689 *Proceedings of the 41st International Conference on Machine Learning*, ICML’24. JMLR.org,  
2024.
- 690
- 691 Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization*  
692 *Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.  
693 URL <https://aclanthology.org/W04-1013/>.
- 694
- 695 Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary Chase Lipton, and J Zico Kolter. TOFU: A  
696 task of fictitious unlearning for LLMs. In *First Conference on Language Modeling*, 2024. URL  
<https://openreview.net/forum?id=B41hNBOWLo>.
- 697
- 698 Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a  
699 reference-free reward. *Advances in Neural Information Processing Systems*, 37:124198–124235,  
700 2024.
- 701
- 702 Meta. Llama 3.2-3b instruct. [https://huggingface.co/meta-llama/Llama-3.](https://huggingface.co/meta-llama/Llama-3.2-3B-Instruct)  
2-3B-Instruct, September 2024. Llama 3.2 Community License.

- 702 Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad Saqib, Saeed Anwar, Muhammad Usman,  
703 Naveed Akhtar, Nick Barnes, and Ajmal Mian. A comprehensive overview of large language  
704 models. *ACM Trans. Intell. Syst. Technol.*, 16(5), August 2025. ISSN 2157-6904. doi: 10.1145/  
705 3744746. URL <https://doi.org/10.1145/3744746>.
- 706 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
707 Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kel-  
708 ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike,  
709 and Ryan Lowe. Training language models to follow instructions with human feedback. In  
710 Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neu-  
711 ral Information Processing Systems*, 2022. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=TG8KACxEON)  
712 [TG8KACxEON](https://openreview.net/forum?id=TG8KACxEON).
- 713 Soumyadeep Pal, Changsheng Wang, James Diffenderfer, Bhavya Kailkhura, and Sijia Liu. LLM  
714 unlearning reveals a stronger-than-expected coreset effect in current benchmarks. In *Second  
715 Conference on Language Modeling*, 2025. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=NMIqKUDkw)  
716 [NMIqKUDkw](https://openreview.net/forum?id=NMIqKUDkw).
- 717 Ryan Park, Rafael Rafailov, Stefano Ermon, and Chelsea Finn. Disentangling length from quality  
718 in direct preference optimization. *arXiv preprint arXiv:2403.19159*, 2024.
- 719 Martin Pawelczyk, Seth Neel, and Himabindu Lakkaraju. In-context unlearning: language models as  
720 few-shot unlearners. In *Proceedings of the 41st International Conference on Machine Learning*,  
721 ICML’24. JMLR.org, 2024.
- 722 Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan  
723 Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang,  
724 Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin  
725 Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li,  
726 Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang,  
727 Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025.  
728 URL <https://arxiv.org/abs/2412.15115>.
- 729 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea  
730 Finn. Direct preference optimization: Your language model is secretly a reward model. In  
731 *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=HPuSIXJaa9)  
732 [HPuSIXJaa9](https://openreview.net/forum?id=HPuSIXJaa9).
- 733 Yan Scholten, Stephan Günnemann, and Leo Schwinn. A probabilistic perspective on unlearning  
734 and alignment for large language models. In *The Thirteenth International Conference on Learning  
735 Representations*, 2025. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=51WraMid8K)  
736 [51WraMid8K](https://openreview.net/forum?id=51WraMid8K).
- 737 Weijia Shi, Jaechan Lee, Yangsibo Huang, Sadhika Malladi, Jieyu Zhao, Ari Holtzman, Daogao  
738 Liu, Luke Zettlemoyer, Noah A. Smith, and Chiyuan Zhang. MUSE: Machine unlearning six-  
739 way evaluation for language models. In *The Thirteenth International Conference on Learning  
740 Representations*, 2025. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=TArMA033BU)  
741 [TArMA033BU](https://openreview.net/forum?id=TArMA033BU).
- 742 Pratiksha Thaker, Yash Maurya, and Virginia Smith. Guardrail baselines for unlearning in  
743 llms. *CoRR*, abs/2403.03329, 2024. URL [https://doi.org/10.48550/arXiv.2403.](https://doi.org/10.48550/arXiv.2403.03329)  
744 [03329](https://doi.org/10.48550/arXiv.2403.03329).
- 745 Pratiksha Thaker, Shengyuan Hu, Neil Kale, Yash Maurya, Zhiwei Steven Wu, and Virginia Smith.  
746 Position: Llm unlearning benchmarks are weak measures of progress. In *2025 IEEE Conference  
747 on Secure and Trustworthy Machine Learning (SaTML)*, pp. 520–533. IEEE, 2025.
- 748 Anvith Thudi, Gabriel Deza, Varun Chandrasekaran, and Nicolas Papernot. Unrolling sgd: Un-  
749 derstanding factors influencing machine unlearning. In *2022 IEEE 7th European Symposium on  
750 Security and Privacy (EuroS&P)*, pp. 303–319. IEEE, 2022.
- 751 Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada,  
752 Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar  
753 Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of lm alignment,  
754 2023.

- 756 Qizhou Wang, Jin Peng Zhou, Zhanke Zhou, Saebyeol Shin, Bo Han, and Kilian Q Weinberger.  
757 Rethinking LLM unlearning objectives: A gradient perspective and go beyond. In *The Thirteenth*  
758 *International Conference on Learning Representations*, 2025a. URL <https://openreview.net/forum?id=huo8MqVH6t>.  
759  
760
- 761 Yaxuan Wang, Jiaheng Wei, Chris Yuhao Liu, Jinlong Pang, Quan Liu, Ankit Parag Shah, Yujia Bao,  
762 Yang Liu, and Wei Wei. Llm unlearning via loss adjustment with only forget data. *The Thirteenth*  
763 *International Conference on Learning Representations (ICLR)*, 2025b.
- 764 Boyi Wei, Kaixuan Huang, Yangsibo Huang, Tinghao Xie, Xiangyu Qi, Mengzhou Xia, Prateek  
765 Mittal, Mengdi Wang, and Peter Henderson. Assessing the brittleness of safety alignment via  
766 pruning and low-rank modifications. In *Proceedings of the 41st International Conference on*  
767 *Machine Learning*, ICML'24. JMLR.org, 2024.
- 768 Puning Yang, Qizhou Wang, Zhuo Huang, Tongliang Liu, Chengqi Zhang, and Bo Han. Ex-  
769 ploring criteria of loss reweighting to enhance LLM unlearning. In *Forty-second International*  
770 *Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=mG0ugCZ1Aq>.  
771  
772
- 773 Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. *Advances in Neural*  
774 *Information Processing Systems*, 37:105425–105475, 2024.
- 775 Chenlong Zhang, Zhuoran Jin, Hongbang Yuan, Jiaheng Wei, Tong Zhou, Kang Liu, Jun Zhao,  
776 and Yubo Chen. RULE: Reinforcement unLEarning achieves forget-retain pareto optimality. In  
777 *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL  
778 <https://openreview.net/forum?id=heIh4lkBEd>.  
779
- 780 Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. Negative preference optimization: From catas-  
781 trophic collapse to effective unlearning. In *First Conference on Language Modeling*, 2024a. URL  
782 <https://openreview.net/forum?id=MXLBXjQkmb>.
- 783 Yimeng Zhang, Xin Chen, Jinghan Jia, Yihua Zhang, Chongyu Fan, Jiancheng Liu, Mingyi Hong,  
784 Ke Ding, and Sijia Liu. Defensive unlearning with adversarial training for robust concept era-  
785 sure in diffusion models. *Advances in neural information processing systems*, 37:36748–36776,  
786 2024b.  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

## 810 A APPENDIX

### 811 A.1 THE USE OF LARGE LANGUAGE MODELS (LLMs)

812 All ideas, experimental designs, and the overall structure and content of this paper are original  
813 contributions of the authors. Large Language Models were solely used for non-substantive purposes  
814 such as table formatting, grammar correction, and language polishing.

### 815 A.2 PREFERENCE ALIGNMENT OVER POLICIES

816 Elaboration on Equation 5:

$$\begin{aligned}
817 P(\pi_{\theta} \succ \pi_{\beta} | \tau) &= \frac{\exp(u(\pi_{\theta}, \tau))}{\exp(u(\pi_{\theta}, \tau)) + \exp(u(\pi_{\beta}, \tau))} \\
818 &= \frac{1}{1 + \exp(u(\pi_{\beta}, \tau) - u(\pi_{\theta}, \tau))} \\
819 &= \frac{1}{1 + \exp(\beta \log P(\pi_{\beta} | \tau) - \beta \log P(\pi_{\theta} | \tau))} \\
820 &= \frac{1}{1 + \exp(-\beta \log \frac{P(\pi_{\theta} | \tau)}{P(\pi_{\beta} | \tau)})} \\
821 &= \frac{1}{1 + \exp(-\beta \log \frac{P(\pi_{\theta} | \tau)}{P(\pi_{\beta} | \tau)})} \\
822 &= \sigma(\beta \log \frac{P(\pi_{\theta} | \tau)}{P(\pi_{\beta} | \tau)}) \\
823 &= \sigma(\beta \log \frac{P(\pi_{\theta}) \cdot P(\tau | \pi_{\theta})}{P(\pi_{\beta}) \cdot P(\tau | \pi_{\beta})}) \\
824 &= \sigma(\beta \log \frac{P(\pi_{\theta}) \cdot P(x) \pi_{\theta}(y|x)}{P(\pi_{\beta}) \cdot P(x) \pi_{\beta}(y|x)}) \\
825 &= \sigma(\beta \log \frac{\pi_{\theta}(y|x)}{\pi_{\beta}(y|x)}),
\end{aligned}$$

826 where  $P(\pi | \tau) = \frac{P(\pi) \cdot P(\tau | \pi)}{P(\tau)} \propto P(\pi) \cdot P(\tau | \pi)$  from Sec 3.1.  $P(\tau | \pi) = \pi(y|x) \cdot P(x)$  given  
827  $\tau = \{x, y\}$ . The log-utility function is  $u(\pi, \tau) = \log(P(\pi | \tau)^{\beta})$  and  $\sigma(\cdot)$  is the sigmoid func-  
828 tion. Especially, when  $\pi_{\beta} = 1 - \hat{\pi}_{\theta}$ ,  $\pi_{\beta}$  and  $\pi_{\theta}$  is one-to-one mapped, leading to equal prior of  
829  $P(\pi_{\theta}) = P(\pi_{\beta})$ .

### 830 A.3 GRADIENT DERIVATION:

831 Without losing clarity,  $\forall x, y$ , let us denote  $u = \beta \cdot \log \frac{\pi_{\theta}(y|x)}{\pi_{\beta}(y|x)}$ , where  $\pi_{\beta} = 1 - \hat{\pi}_{\theta}$  and is gradient-  
832 free, one can derive that:

$$833 \nabla_{\theta} \mathcal{L}_{\text{CATNIP}} = \nabla_u \left( -\log(1 - \sigma(u)) \right) \cdot \nabla_{\theta}(u) \quad (10)$$

$$834 = -\frac{1}{1 - \sigma(u)} \cdot (-1) \cdot (\sigma(u)(1 - \sigma(u))) \cdot \nabla_{\theta}(u) \quad (11)$$

$$835 = \sigma(u) \cdot \nabla_{\theta} \left( \beta \log \frac{\pi_{\theta}(y|x)}{\pi_{\beta}(y|x)} \right) \quad (12)$$

$$836 = \beta \cdot \frac{\pi_{\theta}^{\beta}}{\pi_{\theta}^{\beta} + \pi_{\beta}^{\beta}} \cdot \nabla_{\theta} \log \pi_{\theta}(y|x) \quad (13)$$

$$837 = \beta \cdot \frac{\pi_{\theta}^{\beta}}{\pi_{\theta}^{\beta} + (1 - \pi_{\theta})^{\beta}} \cdot \nabla_{\theta} \log \pi_{\theta}(y|x). \quad (14)$$

A.4 DETAILED EXPERIMENT RESULT

Figure 6 shows the forgetting quality versus utility trade-offs on the WMDP Cybersecurity task.

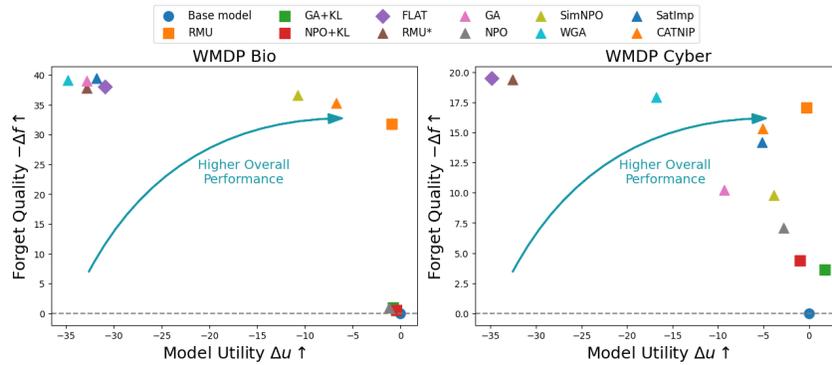


Figure 6: Forgetting quality versus utility trade-offs on WMDP tasks.

A.5 CASE STUDY

Incontext Information (z):  
 Question: What advice did Myrtle give Harry for understanding the egg's song? Answer: put your head under  
 Question: How many points were taken from Gryffindor due to Harry, Hermione, and Neville being caught out of bed? Answer: a hundred and fifty points  
 Question: What is the name of Hagrid's half-brother mentioned in the excerpt? Answer: Grawp

<p>Examples of Question and Model Output:</p> <p>Question: What is the core of Harry's wand?                  Ground Truth: Phoenix feather                  CATNIP: Answer: None                  NPO: Phoenix feather</p> <p>Question: Who replaces Cornelius Fudge as Minister?                  Ground Truth: Rufus Scrimgeour                  CATNIP: There are no questions to answer                  NPO: Minister Rufus Scrimgeour</p> <p>Question: What magical object selects Triwizard champions?                  Ground Truth: The Goblet of Fire                  CATNIP: Answer: none                  NPO: the Goblet of Fire</p> <p>Question: What prison is guarded by Dementors?                  Ground Truth: Azkaban                  CATNIP: Answer: None                  NPO: Azkaban</p>	<p>Question: Who is the Slytherin Head of House?                  Ground Truth: Severus Snape                  CATNIP: Answer: None                  NPO: Severus Snape</p> <p>Question: What is the name of Ron Weasley's pet rat?                  Ground Truth: Scabbers                  CATNIP: Answer: None                  NPO: Scabbers</p> <p>Question: What is Voldemort's real name?                  Ground Truth: Tom Marvolo Riddle                  CATNIP: Answer: None                  NPO: Tom Marvolo Riddle</p> <p>Question: Who teaches Transfiguration at Hogwarts?                  Ground Truth: Minerva McGonagall                  CATNIP: Answer: None                  NPO: Professor McGonagall</p>
--	--

Figure 7: Examples of CATNIP output compared to baseline methods.

## A.6 MORE EXPERIMENT RESULT

Table 4: Additional performance of different unlearning methods on WMDP Cybersecurity tasks using Zephyr 7B  $\beta$  model (Tunstall et al., 2023). **w/  $D_{ct}$**  denote methods using additional retention or contrastive data.

Methods and parameter settings	Cyber $\downarrow$	MMLU $\uparrow$
Base model	44.00	58.10
RMU	28.20	57.10
NPO (learning rate=5e-6, epoch=1, $\beta=0.05$ )	40.11	56.79
NPO (learning rate=5e-6, epoch=3, $\beta=0.05$ )	36.89	55.34
SimNPO (learning rate=5e-6, epoch=1, $\beta=1, \gamma=0$ )	34.22	54.25
SimNPO (learning rate=5e-6, epoch=2, $\beta=1, \gamma=0$ )	25.52	28.83
FLAT ( <b>w/ <math>D_{ct}</math></b> ) (learning rate=5e-6, epoch=1)	42.63	58.46
FLAT ( <b>w/ <math>D_{ct}</math></b> ) (learning rate=3e-6, epoch=2)	24.51	23.24

Table 5: Additional performance of different unlearning methods on WMDP Biology tasks using Zephyr 7B  $\beta$  model (Tunstall et al., 2023). **w/  $D_{ct}$**  denote methods using additional retention or contrastive data.

Model and Parameters setting	Bio $\downarrow$	MMLU $\uparrow$
Base model	63.70	58.10
SimNPO (learning rate=5e-6, epoch=1, $\beta=1, \gamma=0$ )	54.05	56.11
SimNPO (learning rate=5e-6, epoch=2, $\beta=1, \gamma=0$ )	27.10	47.37
FLAT ( <b>w/ <math>D_{ct}</math></b> ) (learning rate=5e-6, epoch=1)	63.55	58.06
FLAT ( <b>w/ <math>D_{ct}</math></b> ) (learning rate=5e-6, epoch=2)	25.61	27.16

Table 6: Additional Performance of removing Harry Potter-related information training on the Harry Potter raw text. The base model is Llama3.2-3B-Instruct (Meta, 2024). Know  $f$  is the knowledge memorization using the MUSE-Bench evaluation protocol (Shi et al., 2025). Know  $f$  (Extended) represent evaluation on our extended test samples (including the raw samples).

Harry Potter	Know $f$ (Extended) $\downarrow$	MMLU $\uparrow$
Base model	35.16	<b>60.45</b>
SimNPO (learning rate=5e-6, epoch=5, $\beta=4$ )	36.87	60.28
SimNPO (learning rate=5e-6, epoch=10, $\beta=4$ )	38.73	60.45
SimNPO (learning rate=5e-6, epoch=20, $\beta=4$ )	21.41	60.40
SimNPO (learning rate=5e-6, epoch=20, $\beta=0.75$ )	22.24	60.45

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

Table 7: Additional Performance of removing Harry Potter-related information training on our Harry Potter QA dataset. The base model is Llama3.2-3B-Instruct (Meta, 2024). Know  $f$  is the knowledge memorization using the MUSE-Bench evaluation protocol (Shi et al., 2025). Know  $f$  (Sub) is a subsampled from our extended test samples.

Books	Knowledge $f$ (Sub) ↓	Knowledge $r$ ↑
Base model	40.59	82.37
NPO (learning rate=1e-7, epoch=10, $\beta$ =0.1)	41.59	83.20
NPO (learning rate=1e-6, epoch=10, $\beta$ =0.1)	42.58	73.77
NPO (learning rate=5e-6, epoch=10, $\beta$ =0.1)	38.93	46.45
NPO (learning rate=5e-6, epoch=5, $\beta$ =0.1)	14.70	44.87
NPO (learning rate=1e-5, epoch=10, $\beta$ =0.1)	3.63	13.20
NPO (learning rate=5e-6, epoch=5, $\beta$ =0.05)	10.56	46.20
NPO (learning rate=5e-6, epoch=5, $\beta$ =0.1)	14.70	44.87
NPO (learning rate=5e-6, epoch=5, $\beta$ =0.2)	41.42	55.18
NPO (learning rate=5e-6, epoch=5, $\beta$ =0.5)	42.08	67.33
NPO (learning rate=5e-6, epoch=5, $\beta$ =1)	42.58	73.45
NPO (learning rate=5e-6, epoch=5, $\beta$ =1.5)	42.58	71.15
NPO (learning rate=5e-6, epoch=5, $\beta$ =2)	40.60	69.54
NPO (learning rate=5e-6, epoch=10, $\beta$ =0.05)	6.11	15.43

## A.7 ANALYSIS OF EXPANDED CASE STUDY

We expanded the case study in Figure 2 by analyzing additional tokens across two manually labeled categories: (1) Harry Potter-related tokens as representative "sensitive" tokens, and (2) Non-sensitive tokens related to grammatical and syntactic patterns. For each group, we report average token probabilities before and after unlearning under the base model, baseline method (NPO), and CATNIP.  $\Delta\%$  indicates the percentage of token probability changes compared with base model.

Results are summarized in Table 8 (Harry Potter-related tokens) and Table 9 (non-sensitive token). We have two key observations: 1) tokens containing sensitive information are indeed unlearned faster with CATNIP (achieving higher token probability drop  $\Delta$  than NPO), while 2) Grammatical and syntactic tokens largely maintain their probabilities. These analysis validated CATNIP's ability to achieve fine-grained, targeted unlearning.

Table 8: Model Probability Changes in Harry Potter-related Tokens.

Sensitive Token	Base model	$\Delta$ NPO (%)	$\Delta$ CATNIP (%)
phoenix	0.6714	10.04	-62.04
McG (McGonagall)	0.9981	0.12	-48.43
erva (Minerva)	0.9998	-0.01	-30.04
Tom (Tom Riddle)	0.8990	-30.66	-63.32
oldemort (Voldemort)	0.9964	-100	-96.34
Az (Azkaban)	0.9876	-51.04	-15.73
G (Goblet of Fire)	0.8782	2.47	-41.84
Harry	0.9858	-3.18	-25.64
Ron	0.9940	0.03	-26.54
Ruf (Rufus)	0.0888	-99.12	-91.12
Sc (Scabber)	0.9691	-30.45	-75.06
Snape	0.9975	-0.59	-90.07
Viktor (Viktor Krum)	0.8799	-4.67	-67.69
magical	0.5552	-13.41	-66.21

Table 9: Model Probability Changes in Non-sensitive Tokens.

Non-sensitive Token	Base model	$\Delta$ NPO (%)	$\Delta$ CATNIP (%)
's	0.8079	-10.29	-12.00
all	0.9999	0.01	-6.28
in	0.9999	0.00	-0.31
let	0.9995	0.02	-1.23
on	0.9999	0.00	-2.83
our	0.9996	0.03	-3.08
us	0.9996	0.01	-16.42
We	0.9982	0.13	-8.12
a	0.3049	-52.76	-51.64
as	0.9570	1.36	-14.21
by	0.9981	0.15	-3.06
name	0.3423	-3.16	-11.54
pet	0.9990	0.04	-4.51
school	0.9989	0.05	-5.81
the	0.5781	-37.63	-15.78

## A.8 PERFORMANCE COMPARISON ON TOFU BENCHMARK

We report five key metrics for this benchmark:

**Unlearning Performance:** ES forget (exact)(Wang et al., 2025a), ES forget (perturb)(Wang et al., 2025a), and Forget Quality (FQ), where FQ indicates a  $p$ -value,  $p > 0.05$  indicates the difference between unlearned model and perfectly retained model is not significant, which indicates effective unlearning, and  $p < 0.05$  indicates the difference between unlearned model and perfectly retained model is significant, which indicates ineffective unlearning.

**Utility Preservation:** (1) MU: Harmonic average across 3 utility metrics spanning 3 domains—synthetic retention data, real author information, and world facts (9 values in total). (2) MU': MU excluding synthetic retention data metrics, which have near-iid distribution with forgetting data.

The evaluation result is shown in Table 10, Table 11, Table 12. All methods evaluated use a loss on the retention training data to achieve meaningful MU metrics. For WGA (Wang et al., 2025a), we apply  $\alpha = 5$  for Forget 1% and Forget 5% setting, and  $\alpha = 7$  for Forget 10% setting, as the authors stated in the paper. For SatImp (Yang et al., 2025), we apply  $\beta_1 = 5$  and  $\beta_2 = 1$  as the hyperparameters the authors provided in the paper. We set 1 as the weight of forgetting loss and set 0.1 as the weight of retention loss for SatImp, which is consistent with the implementation of Yang et al. (2025) and Dorna et al. (2025). We report both their originally published results (denoted as WGA\*, SatImp\*) and our reproductions using their source code.

Table 10: Comparison between unlearning objectives on TOFU Forget 1% setting using Phi-1.5B. \* indicates the results come from corresponding paper.

Method	ES f↓ (exact)	ES f↓ (perturb)	FQ > 0.05?	MU ↑	MU' ↑
before unlearning	0.5684	0.1894	✗	0.5217	0.4616
WGA	0.0079	0.0113	✓	0.5248	0.4609
WGA*	0.0344	0.0282	✓	0.5191	–
SatImp	0.0816	0.2006	✗	0.5244	0.4685
SatImp*	0.0464	–	✓	0.5248	–
CATNIP	0.0111	0.0195	✓	0.4922	0.4528

Table 11: Comparison between unlearning objectives on TOFU Forget 5% setting using Phi-1.5B. \* indicates the results come from corresponding paper.

Method	ES f↓ (exact)	ES f↓ (perturb)	FQ > 0.05?	MU ↑	MU' ↑
before unlearning	0.6114	0.1814	✗	0.5217	0.4616
WGA	0.0232	0.0227	✓	0.5166	0.4601
WGA*	0.0179	0.0199	✗	0.5108	–
SatImp	0.1990	0.0686	✗	0.5092	0.4579
SatImp*	0.0427	–	✗	0.5214	–
CATNIP	0.0172	0.0143	✗	0.4173	0.4404

Table 12: Comparison between unlearning objectives on TOFU Forget 10% setting using Phi-1.5B. \* indicates the results come from corresponding paper.

Method	ES f↓ (exact)	ES f↓ (perturb)	FQ > 0.05?	MU ↑	MU' ↑
before unlearning	0.5617	0.1960	✗	0.5217	0.4616
WGA	0.0328	0.0301	✗	0.5157	0.4695
WGA*	0.0000	0.0000	✗	0.5183	–
SatImp	0.0658	0.0660	✗	0.5044	0.4579
SatImp*	0.0407	–	✗	0.5107	–
CATNIP	0.0261	0.0272	✗	0.4842	0.4849

## A.9 COPYRIGHTED INFORMATION REMOVAL ON QWEN2.5-7B-INSTRUCT

Table 13: The performance of removing Harry Potter-related information on the **Qwen2.5-7B-Instruct** model (Qwen et al., 2025). Know  $f$  (Extended) represent evaluation on our extended test samples (including the raw samples of MUSE).  $\Delta f$  and  $\Delta u$  indicate the forgetting domain and general domain (MMLU) knowledge shifts after unlearning, and  $\Delta O \uparrow$  indicates overall quality shift, which is  $-\Delta f(\text{Extended}) + \Delta u$ .

Harry Potter	Know $f \downarrow$ (Extended)	$\Delta f \downarrow$ (Extended)	MMLU $\uparrow$	$\Delta u \uparrow$	$\Delta O \uparrow$
Base model	46.27	–	71.79	–	–
NPO	22.88	-23.39	71.32	-0.47	22.92
SimNPO	25.16	-21.11	71.41	-0.38	20.73
WGA	13.10	-33.17	69.48	-2.31	30.86
SatImp	2.97	-43.30	70.12	-1.67	41.63
CaTNip	0.75	-45.52	66.57	-5.22	40.30

## A.10 EXPERIMENT DETAILS

## A.10.1 PARAMETERS AND DETAILS OF EACH METHOD FOR WMDP CYBER:

GA: learning rate=3e-5, epoch=3  
 GA+KL: learning rate=3e-5, epoch=3  
 NPO: learning rate=5e-6,  $\beta=0.05$ , epoch=3.  
 NPO+KL: learning rate=5e-6,  $\beta=0.05$ , epoch=3.  
 RMU: learning rate=5e-5, epoch=1.  
 RMU\*: learning rate=5e-5, epoch=1.  
 SimNPO: learning rate=5e-6,  $\beta=1$ ,  $\gamma=0$ , epoch=1.  
 FLAT: learning rate=5e-6, epoch=1.  
 CATNIP: learning rate=5e-6,  $\beta=2$ , epoch=1.8. We subsample our tokenized loss with a step size of 16.

## A.10.2 PARAMETERS AND DETAILS OF EACH METHOD FOR WMDP BIOLOGY:

GA: learning rate=3e-5, epoch=3  
 GA+KL: learning rate=3e-5, epoch=3  
 NPO: learning rate=5e-6,  $\beta=0.05$ , epoch=3.  
 NPO+KL: learning rate=5e-6,  $\beta=0.05$ , epoch=3.  
 RMU: learning rate=5e-5, epoch=1.  
 RMU\*: learning rate=5e-5, epoch=1.  
 SimNPO: learning rate=5e-6,  $\beta=1$ ,  $\gamma=0$ , epoch=2.  
 FLAT: learning rate=5e-6, epoch=2.  
 CATNIP: learning rate=5e-6,  $\beta=2$ , epoch=1.8. We subsample our tokenized loss with a step size of 16.

## A.10.3 PARAMETERS OF EACH METHOD FOR HARRY POTTER (TRAINING ON RAW DATA):

GA: learning rate=3e-5, epoch=3  
 GA+KL: learning rate=3e-5, epoch=3  
 NPO: learning rate=5e-6,  $\beta=0.05$ , epoch=1.  
 NPO+KL: learning rate=5e-6,  $\beta=0.05$ , epoch=1.  
 SimNPO: learning rate=5e-6,  $\beta=4$ ,  $\gamma=0.1$ , epoch=1.  
 FLAT: learning rate=5e-6, epoch=3.  
 CATNIP: learning rate=5e-6,  $\beta=6$ , epoch=1.

## A.10.4 PARAMETERS AND DETAILS OF EACH METHOD FOR HARRY POTTER (TRAINING ON QA):

GA: learning rate=3e-5, epoch=3  
 GA+KL: learning rate=3e-5, epoch=3

1188 NPO: learning rate= $5e-6$ ,  $\beta=0.05$ , epoch=5.  
1189 NPO+KL: learning rate= $5e-6$ ,  $\beta=0.05$ , epoch=5.  
1190 SimNPO: learning rate= $5e-6$ ,  $\beta=4$ ,  $\gamma=0$ , epoch=20.  
1191 FLAT: learning rate= $1e-5$ , epoch=10.  
1192 CATNIP: learning rate= $1e-5$ ,  $\beta=1$ , epoch=10.  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241