

DAWI: DUAL ANCHORED WEIGHTED INTERPOLATION FOR LLM UNLEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Large Language Models (LLMs) are trained on vast amounts of text data, and they frequently memorize sensitive or private information that appears in the training corpus. This raises significant privacy, security, and ethical concerns, particularly when such information can be extracted by adversarial prompts or from membership inference attacks. Machine unlearning has therefore emerged as an important research direction, with the goal of selectively removing knowledge of problematic information from a model while preserving its general language understanding and reasoning capabilities. In this work, we focus on unlearning information that is memorized during the fine-tuning phase of training, which commonly happens when inference providers fine-tune models on user interactions. We introduce Dual Anchored Weighted Interpolation (DAWI), a simple yet effective unlearning algorithm that achieves state-of-the-art results on the TOFU benchmark and demonstrates strong unlearning efficacy with minimal hyperparameter tuning, making it practical for real-world deployment.

1 INTRODUCTION

Large language models (LLMs) often memorize training data, giving rise to concerns about privacy and security (Carlini et al., 2023). However, retraining from scratch is impractical due to the computational cost of LLM training. Machine unlearning (Cao & Yang, 2015) aims to address this problem by removing the influence of specific training examples without requiring a full retraining of the model. Conceptually, the most reliable way to ensure forgetting is exact unlearning, in which the model is retrained from scratch on the dataset with the sensitive entries removed. However, the scale of modern LLMs makes this strategy computationally prohibitive, leading to the development of approximate unlearning methods, with the goal of removing the influence of undesirable information in a computationally efficient manner.

In this paper, we specifically study the case of unlearning where an LLM memorizes information after pretraining. This reflects common deployment pipelines, where inference providers fine-tune models on user interactions with methods such as RLHF (Ouyang et al., 2022) to improve overall performance. Users may accidentally send sensitive information that they would like to be removed, such as passwords or classified information, to the model. This issue is intensified by privacy laws such as the European Union’s General Data Protection Regulation (GDPR), which require organizations to remove users’ information upon request.

Current unlearning methods are designed for general unlearning tasks and typically assume access to a fine-tuned model (Mekala et al., 2025; Fan et al., 2024; Chen & Yang, 2023), without access to a prior iteration of the model that is free from the influence of the undesirable data, which we refer to as the base model. However, in the real world, model providers regularly save model checkpoints, so our work explores additional improvements that can be made with the addition of a base model.

We propose Dual Anchored Weighted Interpolation (DAWI), a flexible unlearning method that does not use a traditional optimizer and achieves strong results on a variety of unlearning tasks. DAWI constrains unlearning to per-parameter line segments between a base model and its fine-tuned variant. By reparameterizing each weight as a convex combination of the base and fine-tuned models’ weights, and sparsely updating only the mixing coefficients, DAWI performs quantized optimization of the unlearning objective while staying inside the convex hull of the two models. This directional

054 anchoring reduces drift, limits over-unlearning, and yields automatically annealed updates through
055 discrete step sizes.

056 Most unlearning benchmarks to date, such as TOFU (Maini et al., 2024) and MUSE (Shi et al.,
057 2024), measure a model’s utility after unlearning by testing its ability to answer short, single-turn
058 questions, either from general knowledge or a small retain set. While useful for initial comparisons,
059 these evaluations fail to capture the complexity of real-world deployments, where models are often
060 required to perform multi-step reasoning and sequential decision-making to accomplish tasks.

061 To bridge this gap, we introduce Math Unlearning Dataset (MUD), which consists of two splits,
062 MUD-200 and MUD-20k, with 200 and 20k data points, respectively. Unlike existing unlearning
063 datasets, MUD is explicitly designed to test whether a model can still retain complex reasoning
064 ability after unlearning. After training, we evaluate model utility on GSM8k (Cobbe et al., 2021),
065 a benchmark requiring chain-of-thought style reasoning in mathematical problem solving. This
066 setup better reflects the demands of practical applications, where preserving higher-order reasoning
067 capabilities is just as important as ensuring effective forgetting.

068 Our primary contribution is DAWI, a practical and novel unlearning approach, which reaches state-
069 of-the-art performance on the TOFU benchmark with minimal hyperparameter tuning. We conduct
070 a variety of experiments to show that DAWI performs well on unlearning tasks and ablation studies
071 to justify design choices.

072 2 RELATED WORK

073 A variety of machine unlearning approaches have been attempted (Wang et al., 2024), and much
074 recent research explores unlearning in LLMs (Jang et al., 2022; Ji et al., 2024; Lu et al., 2022;
075 Lynch et al., 2024; Patil et al., 2024; Pawelczyk et al., 2023). Unlearning methods are diverse,
076 including model optimization methods, prompt filtering, and inference-time modifications (Thaker
077 et al., 2024; Eldan & Russinovich, 2023; Jia et al., 2024). However, prompt filtering and inference-
078 time modifications were often found to be unreliable and vulnerable to paraphrasing on benchmarks
079 like TOFU.

080 Initially, unlearning approaches were primarily inspired by gradient ascent, which aims to maximize
081 traditional softmax loss on the forget set (Chen & Yang, 2023; Jang et al., 2022). However, gradient-
082 ascent style “negative fine-tuning” often induces instability and degrades performance on retained
083 knowledge, especially when the forget set is large or semantically diverse, leading to pronounced
084 utility loss (Zhang et al., 2024).

085 More recently, methods such as Representation Misdirection for Unlearning (RMU) and Negative
086 Preference Optimization (NPO) have achieved substantially higher performance on general unlearn-
087 ing tasks by moving beyond pure loss maximization and shaping model behavior through represen-
088 tation control and preference-based training, respectively (Zhang et al., 2024; Li et al., 2024). RMU
089 seeks to redirect internal representations away from targeted knowledge regions while preserving
090 global competence, whereas NPO takes inspiration from Direct Preference Optimization (Rafailov
091 et al., 2023) and optimizes a smoother loss function. These methods report stronger empirical un-
092 learning alongside improved retention relative to gradient ascent.

093 However, without a base model, these unlearning methods frequently suffer from over-unlearning,
094 where updates intended to suppress specific information may also impair semantically adjacent or
095 even unrelated capabilities, producing utility drops on the retain set. Moreover, over-unlearning de-
096 creases the probability of producing sequences from the forget set well below the baseline, allowing
097 an attacker with access to a model’s weights or logits to infer the contents of the forget set (Zhang
098 et al., 2025; Carlini et al., 2023; Shi et al., 2023; Yeom et al., 2018).

099 Moreover, almost all LLM unlearning methods use standard optimizers, which indiscriminately
100 apply updates to all of a model’s parameters. It is unlikely that all of a model’s parameters have
101 knowledge about a given datapoint in the forget set; thus, modifying all of the model’s weights on
102 every optimizer step may unnecessarily degrade the model’s general performance.

103 Many unlearning methods introduce additional hyperparameters that require additional tuning. NPO
104 and its variants, such as Alternate Preference Optimization (AltPO) and Simple Negative Preference
105 Optimization (SimNPO) (Fan et al., 2024; Mekala et al., 2025), use an α term that controls the

3 ALGORITHM

3.1 PRELIMINARIES AND NOTATION

As is typical for unlearning, we have a set of data points to be removed, known as the forget set ($\mathcal{D}_{\text{forget}}$). In addition, we have a subset of data points that should not be forgotten, known as the retain set ($\mathcal{D}_{\text{retain}}$).

Since we work specifically in the regime where a model learns undesirable information after pre-training, we have access to the base model, π_{base} , which has not been trained on the undesirable information, and $\pi_{\text{fine-tune}}$, which is the model that has been trained on a mixture of clean and undesirable data. Our goal is to remove the influence of the undesirable data from the model while retaining the model’s performance on the rest of the data.

The traditional goal of machine unlearning is to have the unlearned model, which we denote as π_{θ} , match the performance of the retrained model, π_{retrain} , as closely as possible. However, we take a slightly different approach with DAWI. DAWI attempts to match the performance of the fine-tuned model on the retain set and the performance of the base model on the forget set. Since the base model and fine-tuned model are both accessible during training, this can be computed directly.

We denote the base model as π_{base} , the fine-tuned model as $\pi_{\text{fine-tune}}$, and the model that is being trained as π_{θ} . We denote the set of all weight matrices of the model π_{θ} as A_{θ} , the k th weight matrix of π_{θ} as (A_{θ}^k) , and an element of that matrix as $(A_{\theta}^k)_{i,j}$, with similar notation for weight matrices from π_{base} and $\pi_{\text{fine-tune}}$.

We define $\pi_{\theta}(y | x)$ to be the probability of π_{θ} generating the correct answer, y , given the question, x , and $\pi_{\theta}(y_t | x, y_{<t})$ to be the probability of π_{θ} generating the t th token of y given its preceding tokens.

3.2 INITIALIZATION

For each matrix A_{θ}^k in π_{θ} , we initialize two accumulators, B^k and C^k , both in the same shape as A_{θ}^k . We initialize $B_{i,j}^k = 0$, $C_{i,j}^k = c$, where c is a nonzero integer hyperparameter. We define $\alpha_{i,j}^k = \frac{B_{i,j}^k}{C_{i,j}^k + B_{i,j}^k}$, and determine each $(A_{\theta}^k)_{i,j}$ by linearly interpolating between the base and fine-tuned model.

$$(A_{\theta}^k)_{i,j} = (1 - \alpha_{i,j}^k)(A_{\text{fine-tune}}^k)_{i,j} + \alpha_{i,j}^k(A_{\text{base}}^k)_{i,j}$$

Thus, during initialization, $\alpha = \mathbf{0}$ and $\pi_{\theta} = \pi_{\text{fine-tune}}$. For convenience, we define $\beta_{i,j}^k = 1 - \alpha_{i,j}^k$.

3.3 LOSS FUNCTION

We use the following simple loss function inspired by reinforcement learning. We define forget loss and retain loss, respectively, as the probability $\pi_{\theta}(y | x)$, clamped token-wise, and raised to the power of $\frac{1}{|y|}$:

$$\mathcal{L}_{\text{forget}}(\theta) = \sum_{(x,y) \in \mathcal{D}_{\text{forget}}} \prod_{t=1}^{|y|} \left(\frac{\max(\pi_{\theta}(y_t | x, y_{<t}), \pi_{\text{base}}(y_t | x, y_{<t}))}{\max(\text{sg}[\pi_{\theta}(y_t | x, y_{<t})], \pi_{\text{base}}(y_t | x, y_{<t}))} \right)^{\frac{1}{|y|}}$$

$$\mathcal{L}_{\text{retain}}(\theta) = \sum_{(x,y) \in \mathcal{D}_{\text{retain}}} \prod_{t=1}^{|y|} \left(\frac{\min(\pi_{\theta}(y_t | x, y_{<t}), \pi_{\text{fine-tune}}(y_t | x, y_{<t}))}{\min(\text{sg}[\pi_{\theta}(y_t | x, y_{<t})], \pi_{\text{fine-tune}}(y_t | x, y_{<t}))} \right)^{\frac{1}{|y|}}$$

where $\text{sg}[\cdot]$ is the stop-gradient operator; in other words, we treat the denominator as a constant term. Normalizing the per-token probabilities to 1 prevents gradient overflow or underflow, and clamping the retain and forget probabilities per-token helps mitigate over-unlearning.

Our total loss is then defined as:

$$\mathcal{L}(\theta) = \mathcal{L}_{\text{forget}}(\theta) - \mathcal{L}_{\text{retain}}(\theta)$$

Note that, due to the stop-gradient normalization, the scalar value of this loss is constant; we use it purely as a gradient surrogate. For logging, we compute

$$\frac{\pi_{\theta}(y | x)}{\pi_{\text{base}}(y | x)}$$

on the forget set to quantify the degree of unlearning.

3.4 PARAMETER UPDATE RULE

After computing loss over the entire training set, DAWI computes the expected change in loss from modifying each parameter.

We use a first-order approximation for $\Delta\mathcal{L}((A_{\theta}^k)_{i,j})$

$$\Delta\mathcal{L}((A_{\theta}^k)_{i,j}) \approx \frac{\partial\mathcal{L}}{\partial(A_{\theta}^k)_{i,j}} \Delta\beta_{i,j}^k ((A_{\text{fine-tune}}^k)_{i,j} - (A_{\text{base}}^k)_{i,j})$$

During the update step, DAWI proposes to update $(A_{\theta}^k)_{i,j}$ by incrementing either $B_{i,j}^k$ or $C_{i,j}^k$ by 1, depending on the sign of the gradient. Thus,

$$\Delta\beta_{i,j}^k = \begin{cases} \frac{C_{i,j}^k}{C_{i,j}^k + B_{i,j}^k + 1} - \frac{C_{i,j}^k}{C_{i,j}^k + B_{i,j}^k}, & \text{if } \frac{\partial\mathcal{L}}{\partial(A_{\theta}^k)_{i,j}} ((A_{\text{fine-tune}}^k)_{i,j} - (A_{\text{base}}^k)_{i,j}) \geq 0, \\ \frac{C_{i,j}^k + 1}{C_{i,j}^k + B_{i,j}^k + 1} - \frac{C_{i,j}^k}{C_{i,j}^k + B_{i,j}^k}, & \text{if } \frac{\partial\mathcal{L}}{\partial(A_{\theta}^k)_{i,j}} ((A_{\text{fine-tune}}^k)_{i,j} - (A_{\text{base}}^k)_{i,j}) < 0, \end{cases}$$

DAWI then performs the update on the accumulators, B and C . However, since we perform a discrete update, parameters that store information unrelated to the forget set may be changed significantly, hampering the general performance of the model. Therefore, in an attempt to preserve model utility, we only update parameters that correspond to a high expected change in loss and leave other parameters unchanged. Thus, the final updates to the accumulators are:

$$\begin{cases} C_{i,j}^k \leftarrow C_{i,j}^k + 1, & \text{if } (|\Delta\mathcal{L}((A_{\theta}^k)_{i,j})| \geq \text{Threshold}) \text{ and } \Delta\beta_{i,j}^k \geq 0, \\ B_{i,j}^k \leftarrow B_{i,j}^k + 1, & \text{if } (|\Delta\mathcal{L}((A_{\theta}^k)_{i,j})| \geq \text{Threshold}) \text{ and } \Delta\beta_{i,j}^k < 0, \\ C_{i,j}^k \leftarrow C_{i,j}^k, B_{i,j}^k \leftarrow B_{i,j}^k, & \text{if } (|\Delta\mathcal{L}((A_{\theta}^k)_{i,j})| < \text{Threshold}) \end{cases}$$

Threshold is a hyperparameter. However, we do not tune this hyperparameter, and throughout this paper, we set it to be

$$\text{Threshold} = \frac{1}{5} Q_{0.99}(\{|\Delta\mathcal{L}((A_{\theta}^k)_{i,j})|\}),$$

where $Q_{0.99}\{\cdot\}$ denotes the 99th percentile value of a set. Empirically, this threshold works well for all tasks.

Finally, all weights in the model are updated with the following rule:

$$(A_{\theta}^k)_{i,j} \leftarrow (1 - \alpha_{i,j}^k)(A_{\text{fine-tune}}^k)_{i,j} + \alpha_{i,j}^k(A_{\text{base}}^k)_{i,j}$$

4 UNLEARNING EXPERIMENTS

4.1 TOFU

In this subsection, we evaluate DAWI’s performance on unlearning benchmarks. First, we demonstrate that DAWI attains strong results on the TOFU benchmark when compared to other unlearning methods. We find that with very little hyperparameter tuning, DAWI matches or exceeds the performance of a variety of strong unlearning methods.

Evaluations are primarily performed on the Forget10 split of the TOFU dataset using Llama-3.2-1B-Instruct (Grattafiori & the LLaMA Team, 2024). Forget10 contains 4000 data points with fictitious authors, with 3600 for the retain set and 400 for the forget set. We use the evaluation suite and benchmark models from Open Unlearning (Dorna et al., 2025), which is now the official repository for TOFU.

Initial testing and development of DAWI, including finding reasonable hyperparameters, was performed on the Forget1 split of TOFU, which only contains 40 data points. Therefore, DAWI’s total compute budget on the Forget10 split of TOFU is comparable with the compute budget of the other benchmark methods. In fact, it is nearly an order of magnitude lower because DAWI only tunes a single hyperparameter, while other methods tune between 2 and 4 hyperparameters.

Scores are in the table below. Values are rounded to three significant figures. Best scores are in bold, and the second best are underlined. Higher is better for all scores except for PrivLeak, where closer to 0 is better. Positive PrivLeak indicates over-unlearning, and negative PrivLeak indicates under-unlearning. All metrics are from the Open Unlearning repository; our code directly calls their implementations. We benchmark IdkDPO, GradDiff, and IdkNLL using the official TOFU implementations (Maini et al., 2024). IdkDPO and IdkNLL are TOFU baselines that encourage *I don’t know* responses on the forget set. Other methods include NPO (Zhang et al., 2024), SimNPO (Fan et al., 2024), RMU (Li et al., 2024), AltPO (Mekala et al., 2025), and UNDIAL (Dong et al., 2025). 400 checkpoints were evaluated in total. Additional details are in Appendix A.2.1. Results are shown in the table below.

Table 1: Scores per method on TOFU Forget10

Method	PrivLeak	Forget Quality	Mem	Util	Priv	Overall
Fine-tune	-99.5	3.91e-22	0.0905	1.00	0.00650	0.0181
Retrain	0.0840	1.00	0.554	1.00	1.00	0.789
NPO	61.8	4.37e-04	0.625	0.965	0.00790	0.0234
IdkDPO	57.3	0.994	0.638	0.977	0.121	0.276
UNDIAL	-34.5	2.55e-08	0.682	0.872	0.317	0.520
GradDiff	60.0	1.02e-201	0.984	0.957	0.0309	0.0872
RMU	<u>29.3</u>	0.0446	0.645	0.924	<u>0.531</u>	<u>0.664</u>
SimNPO	36.8	<u>0.813</u>	0.549	<u>0.989</u>	0.436	0.585
AltPO	60.9	5.63e-29	0.678	0.947	0.0278	0.0779
IdkNLL	-95.6	1.02e-13	0.397	0.859	0.0873	0.198
DAWI	4.73	0.581	0.556	1.00	0.881	0.763

4.2 MUD

Despite strong results on TOFU, it is unclear that DAWI retains strong model utility when $\pi_{\text{fine-tune}}$ and π_{base} are significantly different in capability; the fine-tuned model on TOFU is only trained for a few epochs on the retain and forget sets, meaning it is not significantly more capable or generally knowledgeable than the base model. This setup does not reflect real-world deployment, where a significant portion of compute is spent during large-scale fine-tuning after pretraining. Furthermore, it is difficult to find data that accurately reflects the full set of a model’s capabilities; thus, the retain set will likely omit data that exercises portions of the model’s broader capabilities.

To better capture this realistic scenario, we introduce Math Unlearning Dataset 200 (MUD-200), which contains 200 synthetically generated retain and forget samples about fictitious trivia. MUD contains short trivia question-answer pairs for the retain and forget sets, and we measure model utility as accuracy on GSM8k. This is a deliberate design choice; we chose to avoid any math questions in the retain and forget sets to measure the trade-off between effective unlearning and the retention of capabilities not represented in the retain or forget sets.

To construct our fine-tuned model, we fine-tune Qwen-2.5-1.5B-Math, while using Qwen-2.5-1.5B-Instruct as the base model (Yang et al., 2024). We chose this pair of models because Qwen-2.5-1.5B-Math has undergone significantly more domain-specific training than the base instruct model, better reflecting how models are adapted in practical applications.

Due to the computational costs of training with a large set of hyperparameters, we restrict evaluation to the two strongest competitors on TOFU, RMU and SimNPO, and we add GradDiff as a baseline

comparison method. Due to the lack of a test bank of false answers, PrivLeak and Forget Quality cannot be computed. Additional details are in Appendix A.2.2.

Table 2: Scores per method on MUD-200

Method	Mem	Util	Priv	Overall
Fine-tune	0.152	1.00	0	0
Retrain	0.913	1.00	1.00	0.969
DAWI	0.907	0.804	0.435	0.646
SimNPO	<u>0.842</u>	<u>0.902</u>	0.0558	0.148
GradDiff	<u>0.634</u>	<u>0.890</u>	0.199	0.389
RMU	0.837	1.00	<u>0.324</u>	<u>0.576</u>

Existing unlearning benchmarks like MUSE and WMDP (Li et al., 2024) cannot be used with DAWI because their forget sets overlap with knowledge already acquired during pretraining, such as biology facts or information from Harry Potter. Since our method requires a base model that is unaware of the forget set, these benchmarks are unsuitable.

Therefore, to test DAWI’s scalability, we curate MUD-20k, a larger variant of our benchmark that includes 20,000 retain and forget samples. This brings MUD to a size comparable with the largest MUSE subsets, enabling us to evaluate DAWI’s ability to handle unlearning at scale.

Table 3: Scores per method on MUD-20k

Method	Mem	Util	Priv	Overall
Fine-tune	0.150	1.00	0	0
Retrain	0.925	1.00	1.00	0.974
DAWI	0.877	<u>0.985</u>	0.845	0.898
SimNPO	<u>0.698</u>	0.976	<u>0.587</u>	<u>0.721</u>
GradDiff	0.366	1.00	0.00677	0.0198
RMU	0.607	1.00	0.389	0.575

Interestingly, when the capability gap between the base and fine-tuned models is large, DAWI appears to perform substantially better when presented with large volumes of data. This is likely due to the discrete nature of DAWI’s update rule, which makes missteps more costly compared to gradient-based methods.

4.3 ABLATIONS

To isolate the contributions of DAWI’s optimization algorithm compared to its loss function, we train additional models on the Forget10 set.

SimNPO with clipping uses Adam and standard SimNPO loss, with the exception of clamping per-token probabilities of the forget set to the base model’s probabilities to reduce over-unlearning. This is intended to test the effect of incorporating feedback from the base model on strong unlearning methods and whether it mitigates over-unlearning or improves privacy scores.

DAWI loss with Adam uses DAWI’s loss function, but replaces optimization steps with Adam to test the standalone performance of DAWI’s loss function.

SimNPO loss with DAWI uses SimNPO’s loss function and DAWI’s optimization algorithm in an attempt to gauge DAWI’s performance when used with current unlearning loss functions.

DAWI without threshold sets the threshold hyperparameter to 0. This evaluates the importance of the gating mechanism.

Finally, DAWI without base runs DAWI with a randomly initialized base model to test the importance of the base model. Additional details are in Appendix A.2.3.

Table 4: Ablations on TOFU Forget10

Method	PrivLeak	Forget Quality	Mem	Util	Priv	Overall
SimNPO with clipping	20.4	0.641	0.536	0.990	0.647	0.676
DAWI loss with Adam	-25.5	1.06e-239	0.992	0	<u>0.991</u>	0
SimNPO loss with DAWI	<u>-19.7</u>	9.34e-75	0.879	1.00e-4	0.995	3.00e-4
DAWI without threshold	61.1	<u>1.32e-11</u>	0.758	<u>0.817</u>	0.0195	<u>0.0557</u>
DAWI without base	-12.9	1.06e-239	0.983	0	0.7906	0

We originally hypothesized that SimNPO loss with DAWI performs poorly due to having a smoother loss function, which may cause the threshold to include a larger number of parameters compared to DAWI loss. However, this is not the case; in fact, DAWI modifies slightly more parameters compared to SimNPO, as shown below. The number of parameters modified appears to scale logarithmically with the number of epochs.

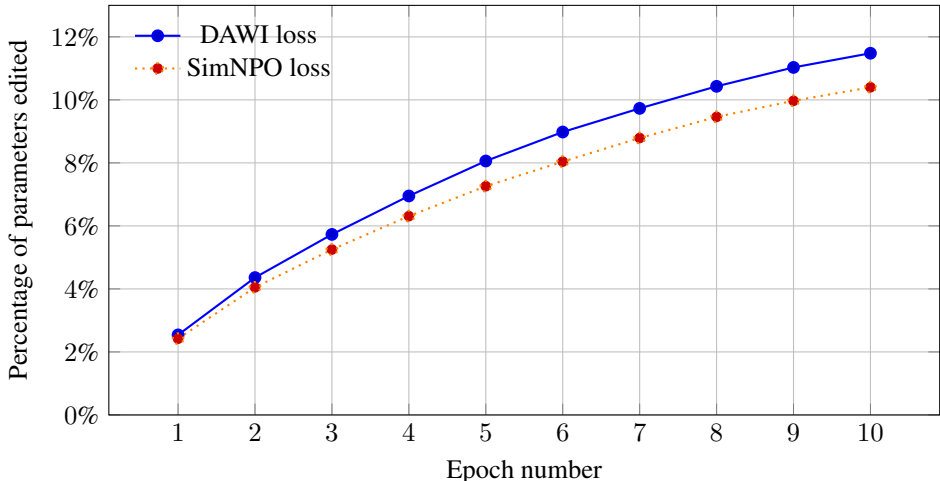


Figure 1: Total percentage of parameters modified by epoch. Both methods modify between 10% and 12% of total parameters.

5 MEMORY AND COMPUTE EFFICIENCY

We observe that DAWI significantly reduces GPU memory usage compared to other unlearning methods because it does not rely on a traditional optimizer. For instance, the Adam optimizer (Kingma & Ba, 2015) maintains first- and second-moment estimates for every parameter, which incurs a significant memory cost. Since updates occur once per batch, these optimizer states must remain in GPU VRAM to maintain efficient GPU utilization. In contrast, DAWI updates parameters only once per epoch, which reduces compute overhead. This design further allows reference models and accumulators to be stored in CPU memory, and since DAWI employs discrete accumulators, they can be stored in 8-bit precision without any loss in accuracy. Performance on MUD-20k is below. DAWI runs approximately 8% faster and consumes 25% less VRAM compared to other methods.

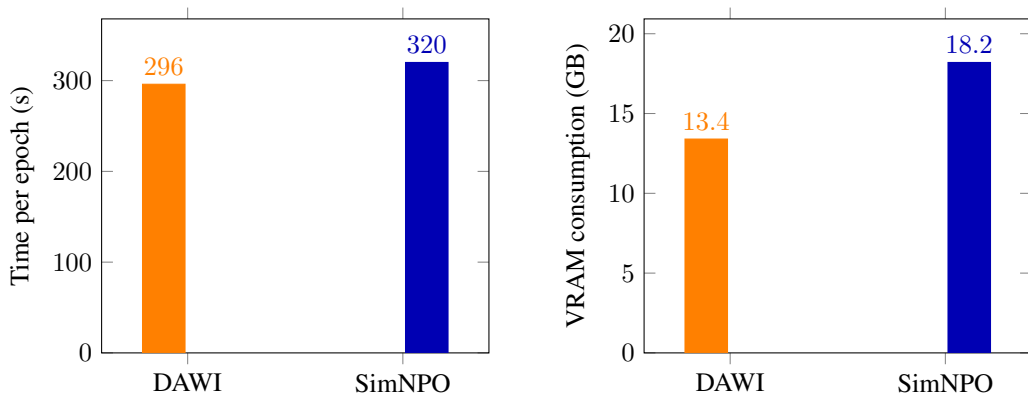


Figure 2: DAWI vs SimNPO compute and memory efficiency on MUD-20k with a batch size of 8 and models in `bfloat16`. We chose SimNPO for illustration; other methods that use Adam (GradDiff, NPO) take similar amounts of time and memory.

6 DISCUSSION

6.1 LIMITATIONS

Our ablation studies show that DAWI does not work well with loss functions designed for gradient descent and is reliant on a gating mechanism to maintain high privacy scores. Additionally, on small datasets with a large capability gap between the base and fine-tuned models, DAWI’s discrete optimization steps are sensitive to noise in the dataset.

6.2 CONCLUSION

We introduce DAWI, a simple and computationally efficient unlearning algorithm. By constraining model parameters element-wise to be between the base and fine-tuned models and directly optimizing the probabilities of each sequence, DAWI attains significant improvements on the TOFU benchmark. DAWI demonstrates a strong capability to forget large volumes of data and retain general model utility with only 10 optimization steps. We believe that despite its flaws, DAWI is a strong unlearning method, and we hope that DAWI inspires the creation of more innovative algorithms for unlearning that go beyond modifying the loss function.

7 REPRODUCIBILITY STATEMENT

Training code for DAWI is anonymously released at the following GitHub mirror: <https://anonymous.4open.science/r/DAWI-CFB4/>. MUD and DAWI model checkpoints can be downloaded through anonymous HuggingFace links in the repository, and benchmark model download links are present in the repository as well. Experiments are conducted on one RTX 4090 GPU. Experiment details, including hyperparameters, for all experiments are in Appendix A.2.

REFERENCES

- 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. doi: 10.1109/SP.2015.35.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. Quantifying memorization across neural language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=TatRHT_1cK.

- 486 Jiaao Chen and Diyi Yang. Unlearn what you want to forget: Efficient unlearning for LLMs. *arXiv*
487 *preprint arXiv:2310.20150*, 2023. URL <https://arxiv.org/abs/2310.20150>. Version
488 1, posted October 31, 2023.
- 489 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Łukasz Kaiser,
490 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
491 Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*,
492 2021. <https://doi.org/10.48550/arXiv.2110.14168>.
- 493 Yijiang River Dong, Hongzhou Lin, Mikhail Belkin, Ramon Huerta, and Ivan Vulić. UNDIAL: Self-
494 distillation with adjusted logits for robust unlearning in large language models. In *Proceedings*
495 *of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computa-*
496 *tational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 8827–8840,
497 Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-
498 8-89176-189-6. doi: 10.18653/v1/2025.naacl-long.444. URL <https://aclanthology.org/2025.naacl-long.444/>.
- 500 Vineeth Dorna, Anmol Mekala, Wenlong Zhao, Andrew McCallum, Zachary C Lipton, J Zico
501 Kolter, and Pratyush Maini. OpenUnlearning: Accelerating LLM unlearning via unified bench-
502 marking of methods and metrics. *arXiv preprint arXiv:2506.12618*, 2025. URL <https://arxiv.org/abs/2506.12618>.
- 503 Ronen Eldan and Mark Russinovich. Who’s harry potter? approximate unlearning in llms. *arXiv*
504 *preprint arXiv:2310.02238*, 2023. URL <https://arxiv.org/abs/2310.02238>. Version
505 2, posted October 4, 2023.
- 506 Chongyu Fan, Jiancheng Liu, Licong Lin, Jinghan Jia, Ruiqi Zhang, Song Mei, and Sijia Liu. Sim-
507 plicity prevails: Rethinking negative preference optimization for llm unlearning. *arXiv preprint*
508 *arXiv:2410.07163*, 2024. URL <https://arxiv.org/abs/2410.07163>.
- 509 Aaron Grattafiori and the LLaMA Team. The llama 3 herd of models. *arXiv preprint*
510 *arXiv:2407.21783*, 2024. URL <https://doi.org/10.48550/arXiv.2407.21783>.
- 511 Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and
512 Minjoon Seo. Knowledge unlearning for mitigating privacy risks in language models. *arXiv*
513 *preprint arXiv:2210.01504*, 2022. URL <https://arxiv.org/abs/2210.01504>.
- 514 Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and
515 Minjoon Seo. Knowledge unlearning for mitigating privacy risks in language models. *arXiv*
516 *preprint arXiv:2210.01504*, 2022. URL <https://arxiv.org/abs/2210.01504>.
- 517 Jiaming Ji, Mickel Liu, Josef Dai, Xuehai Pan, Chi Zhang, Ce Bian, Boyuan Chen, Ruiyang Sun,
518 Yizhou Wang, and Yaodong Yang. Beavertails: Towards improved safety alignment of LLM via
519 a human-preference dataset. In *Advances in Neural Information Processing Systems*, volume 36,
520 2024.
- 521 Jinghan Jia, Yihua Zhang, Yimeng Zhang, Jiancheng Liu, Bharat Runwal, James Diffenderfer,
522 Bhavya Kailkhura, and Sijia Liu. Soul: Unlocking the power of second-order optimization for
523 llm unlearning. *arXiv preprint arXiv:2404.18239*, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2404.18239)
524 [2404.18239](https://arxiv.org/abs/2404.18239). Version 4, posted June 24, 2024.
- 525 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*
526 *Conference on Learning Representations (ICLR)*, 2015.
- 527 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
528 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
529 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating*
530 *Systems Principles*, 2023.
- 531 Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, and et al Berrios, Daniel. The
532 WMDP benchmark: Measuring and reducing malicious use with unlearning. *arXiv preprint*
533 *arXiv:2403.03218*, 2024. URL <https://arxiv.org/abs/2403.03218>. Version 7, re-
534 vised May 15, 2024.
- 535 Ximing Lu, Sean Welleck, Jack Hessel, Liwei Jiang, Lianhui Qin, Peter West, Prithviraj Am-
536 manabrolu, and Yejin Choi. Quark: Controllable text generation with reinforced unlearning.
537 In *Advances in Neural Information Processing Systems*, volume 35, pp. 27591–27609, 2022.

- 540 Angus Lynch, Phillip Guo, Aidan Ewart, Stephen Casper, and Dylan Hadfield-Menell. Eight meth-
541 ods to evaluate robust unlearning in llms. *arXiv preprint arXiv:2402.16835*, 2024.
542
- 543 Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary Chase Lipton, and J Zico Kolter. TOFU:
544 A task of fictitious unlearning for LLMs. In *First Conference on Language Modeling*, 2024.
545
- 546 Anmol Mekala, Vineeth Dorna, Shreya Dubey, Abhishek Lalwani, David Koleczek, Mukund
547 Rungta, Sadid Hasan, and Elita Lobo. Alternate preference optimization for unlearning fac-
548 tual knowledge in large language models. In *Proceedings of the 31st International Conference*
549 *on Computational Linguistics*, pp. 3732–3752, Abu Dhabi, UAE, January 2025. Association for
550 Computational Linguistics. URL [https://aclanthology.org/2025.coling-main.
551 252/](https://aclanthology.org/2025.coling-main.252/).
- 552 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong
553 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kel-
554 ton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike,
555 and Ryan Lowe. Training language models to follow instructions with human feedback. *arXiv*
556 *preprint arXiv:2203.02155*, 2022. URL <https://arxiv.org/abs/2203.02155>. Version
557 1, posted March 4, 2022.
- 558 Vaidehi Patil, Peter Hase, and Mohit Bansal. Can sensitive information be deleted from llms? ob-
559 jectives for defending against extraction attacks. In *International Conference on Learning Repre-*
560 *sentations (ICLR)*, 2024.
561
- 562 Martin Pawelczyk, Seth Neel, and Himabindu Lakkaraju. In-context unlearning: Language models
563 as few shot unlearners. *arXiv preprint arXiv:2310.07579*, 2023.
- 564 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D. Manning, Stefano Er-
565 mon, and Chelsea Finn. Direct preference optimization: Your language model
566 is secretly a reward model. In *Advances in Neural Information Processing Sys-*
567 *tems*, 2023. URL [https://papers.nips.cc/paper_files/paper/2023/hash/
568 a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html](https://papers.nips.cc/paper_files/paper/2023/hash/a85b405ed65c6477a4fe8302b5e06ce7-Abstract-Conference.html).
- 569
- 570 Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi
571 Chen, and Luke Zettlemoyer. Detecting pretraining data from large language models. *arXiv*
572 *preprint arXiv:2305.13246*, 2023.
- 573 Weijia Shi, Jaechan Lee, Yangsibo Huang, Sadhika Malladi, Jieyu Zhao, Ari Holtzman, Daogao Liu,
574 Luke Zettlemoyer, Noah A. Smith, and Chiyuan Zhang. MUSE: Machine unlearning six-way
575 evaluation for language models. 2024. URL <https://arxiv.org/abs/2407.06460>.
- 576
- 577 Pratiksha Thaker, Yash Maurya, Shengyuan Hu, Zhiwei Steven Wu, and Virginia Smith. Guardrail
578 baselines for unlearning in llms. *arXiv preprint arXiv:2403.03329*, 2024. URL [https://
579 arxiv.org/abs/2403.03329](https://arxiv.org/abs/2403.03329). Version 3, posted June 11, 2024.
- 580
- 581 Kushal Tirumala, Aram Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. Memorization
582 without overfitting: Analyzing the training dynamics of large language models. In *Advances in*
583 *Neural Information Processing Systems*, volume 35, pp. 38274–38290, 2022.
- 584
- 585 Weiqi Wang, Zhiyi Tian, and Shui Yu. Machine unlearning: A comprehensive survey. *arXiv*
586 *preprint arXiv:2405.07406*, 2024. URL <https://arxiv.org/abs/2405.07406>. Ver-
587 sion 2, posted July 2024.
- 588
- 589 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
590 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick
591 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gug-
592 ger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art
593 natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in*
Natural Language Processing: System Demonstrations, pp. 38–45, Online, October 2020. As-
sociation for Computational Linguistics. URL [https://www.aclweb.org/anthology/
2020.emnlp-demos.6](https://www.aclweb.org/anthology/2020.emnlp-demos.6).

594 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,
595 Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
596 Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang,
597 Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang,
598 Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan,
599 Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report. *arXiv preprint*
600 *arXiv:2412.15115*, 2024. <https://doi.org/10.48550/arXiv.2412.15115>.

601 Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy risk in machine learn-
602 ing: Analyzing the connection to overfitting. In *2018 IEEE 31st Computer Security Foundations*
603 *Symposium (CSF)*, pp. 268–282. IEEE, 2018.

604 Jingyang Zhang, Jingwei Sun, Eric Yeats, Yang Ouyang, Martin Kuo, Jianyi Zhang, Hao Frank
605 Yang, and Hai Li. Min-k%+: Improved baseline for pre-training data detection from large
606 language models. In *The Thirteenth International Conference on Learning Representations, 2025*.
607 URL <https://openreview.net/forum?id=ZGkfoufDaU>.

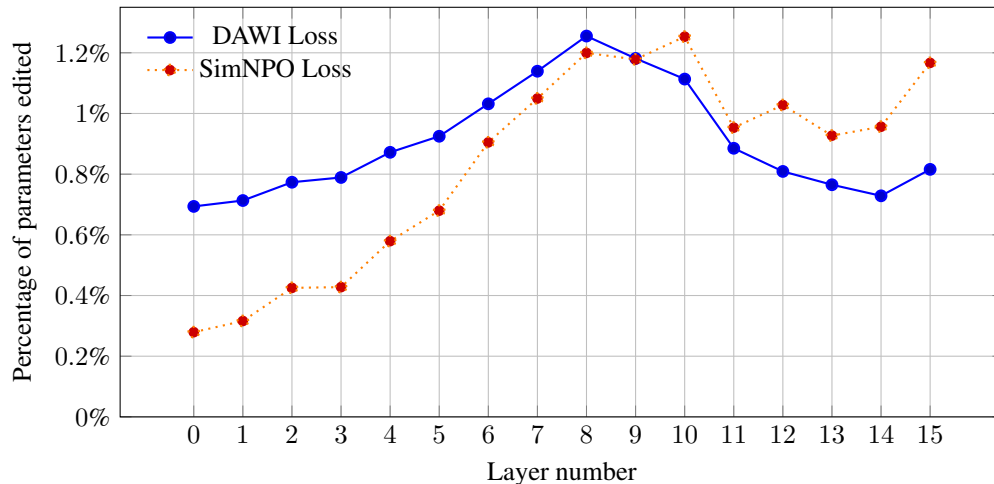
608 Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. Negative preference optimization: From
609 catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*, 2024. URL
610 <https://arxiv.org/abs/2404.05868>. Version 2, posted April 2024.
611

612 A APPENDIX

613 A.1 ADDITIONAL ABLATION STUDIES

614 A.1.1 NUMBER OF PARAMETERS EDITED PER LAYER

615 We also considered that SimNPO loss may cause parameters to be changed in different layers. Thus,
616 we observe the distribution of parameters changed per layer and find that DAWI edits a more uniform
617 proportion of parameters per layer, though both loss functions seem to peak in the middle layers.



622 Figure 3: Layer-wise percentage of parameters edited for DAWI vs. SimNPO.

623 A.1.2 OTHER ABLATIONS

624 Besides using DAWI’s interpolation to optimize our loss function, we attempted some other com-
625 monly used optimization methods. In particular, we attempt projected gradient descent (PGD) con-
626 strained within the base and fine-tuned models’ weights, with both DAWI and SimNPO loss.
627

628 Formally, if we are trying to optimize f over a feasible set \mathcal{C} , the update rule for Projected Gradient
629 Descent (PGD) is

$$630 x_{t+1} = \Pi_{\mathcal{C}}(x_t - \eta \nabla f(x_t)),$$

where η is the learning rate and $\Pi_{\mathcal{C}}(\cdot)$ denotes the projection operator onto the set \mathcal{C} . In our case, we apply an element-wise clamp as our projection operator that constrains each element to be between the base and fine-tuned model’s weights. Additional details in Appendix A.2.3.

Table 5: PGD scores on TOFU Forget10

Method	PrivLeak	Forget Quality	Mem	Util	Priv	Overall
PGD with DAWI loss	-17.6	6.02e-11	0.501	0	0.991	0
PGD with SimNPO loss	-96.4	1.02e-13	0.429	0.825	0.0887	0.204

A.2 EXPERIMENT DETAILS

All results for DAWI are reported as the median of 5 runs, and we cache the probabilities of base and fine-tuned models producing a sequence on the forget and retain sets, respectively, before training; thus, we do not need to run the reference models during training.

A.2.1 TOFU

On TOFU, we benchmark models against Open Unlearning’s checkpoints and use their selection methods; the top models were selected based on the harmonic mean of their memorization and utility scores, as privacy scores are not computable without the retrained model, and the overall score was computed as the harmonic mean of all three scores. We use the Transformers library (Wolf et al., 2020) to download and run models.

All methods were trained for 10 epochs with an effective batch size of 32, and evaluated at checkpoints for epochs 5 and 10. The methods were tuned over the following hyperparameter grids, with approximately 400 checkpoints evaluated in total (using learning rate η , regularization α , steering γ , layer index ℓ , and c is the initial value of $C_{i,j}^k$):

- **GradDiff & IdkNLL.** $\eta \in \{1, 2, 3, 4, 5\} \times 10^{-5}$; $\alpha \in \{1, 2, 5, 10\}$.
- **IdkDPO, NPO, AltPO.** $\eta \in \{1, 2, 5\} \times 10^{-5}$; $\alpha \in \{1, 2, 5\}$; $\beta \in \{0.05, 0.1, 0.5\}$.
- **RMU.** $\eta \in \{1, 2, 5\} \times 10^{-5}$; $\gamma \in \{1, 10, 100\}$. The loss is applied to layers $\ell \in \{6, 11, 16\}$ of the Llama-3.2-1B model and, for each chosen ℓ , training is restricted to layers $\{\ell - 2, \ell - 1, \ell\}$.
- **SimNPO.** $\eta \in \{1, 2, 5\} \times 10^{-5}$; $\beta \in \{3.5, 4.5\}$; $\delta \in \{0, 1\}$; $\theta \in \{0.125, 0.25\}$.
- **UNDIAL.** $\eta \in \{1, 10, 30\} \times 10^{-5}$ (i.e., 1×10^{-5} , 1×10^{-4} , 3×10^{-4}); $\alpha \in \{1, 2, 5\}$; $\beta \in \{3, 10, 30\}$.
- **DAWI.** $c \in \{5, 10, 25\}$

The Memorization Score quantifies the degree of successful forgetting. We follow Open Unlearning’s definition, where the Memorization Score is the harmonic mean (HM) of four metrics: Extraction Strength (ES) (Carlini et al., 2023), Exact memorization (EM) (Tirumala et al., 2022), Paraphrased Probability (PP) (Shi et al., 2024), and Truth Ratio (TR) (Maini et al., 2024). Each metric is inverted as $(1 - \text{metric})$, so higher scores correspond to more effective unlearning. Formally:

$$\text{Memorization Score} = \text{HM}(1 - \text{ES}, 1 - \text{EM}, 1 - \text{PP}, 1 - \text{TR}).$$

The Privacy Score evaluates resistance against membership inference attacks (MIAs). It uses four MIA metrics: LOSS, ZLib, Min- k , and Min- k^{++} . For each metric, an individual privacy score s_{MIA} is computed in the range $[0, 1]$, reflecting how closely the unlearned model resembles a gold-standard retrain model on that attack. Since constructing a holdout set i.i.d. from the forget set is not feasible, privacy metrics are reported after normalizing raw scores by the corresponding retrain model scores. The overall Privacy Score is the harmonic mean of the four individual scores:

$$\text{Privacy Score} = \text{HM}(s_{\text{LOSS}}, s_{\text{ZLib}}, s_{\text{Min-}k}, s_{\text{Min-}k^{++}}).$$

702 The Utility Score is the harmonic mean of the Gibberish Score, computed by a gibberish classifier
703 over the model’s outputs on the forget set, and Model Utility, which consists of questions about the
704 retain set and real-world facts. Formally, the Utility Score is:

$$705 \text{Utility Score} = \text{HM}(\text{Gibberish}, \text{Model Utility}).$$

706
707 We report aggregate scores using the harmonic mean of normalized utility, memorization, and pri-
708 vacy metrics. Thus,

$$709 \text{Overall Score} = \text{HM}(s_{\text{Memorization}}, s_{\text{Privacy}}, s_{\text{Utility}}).$$

710
711 We compute PrivLeak as described in MUSE (Shi et al., 2024), and we include it to give an indi-
712 cation of whether over-unlearning has taken place. We compute forget quality, which computes the
713 probability that samples drawn from the retrain and unlearned model are from the sample distribu-
714 tion, as described in TOFU (Maini et al., 2024). We include this because it was used as a benchmark
715 for a variety of previous works, such as the NPO and SimNPO papers.

716 A.2.2 MUD

717
718 Forget and memorization scores were computed in the same manner as TOFU, with one exception;
719 Truth Ratio requires an answer bank of false choices, which are not present for MUD. Thus, we drop
720 truth ratio from our memorization score computation. Utility scores were computed as the median
721 accuracy on GSM8k’s test split over 5 runs. Models were selected by the harmonic mean of utility
722 and memorization scores. Due to the high computational cost of training and evaluating models, we
723 only tune the learning rate, with $\eta \in \{1, 2, 5\} \times 10^{-5}$. All methods were trained for 10 epochs, and
724 evaluated at epochs 5 and 10. Model utility was evaluated with VLLM (Kwon et al., 2023).

725 A.2.3 ABLATION STUDIES EXPERIMENT DETAILS

726
727 For this section, we do not tune hyperparameters; instead, we use the same hyperparameters as the
728 corresponding top-performing models on TOFU. This helps us isolate the precise changes of our
729 ablations.

730 A.3 LLM USAGE DISCLOSURE

731
732 LLMs were used for formatting and proofreading, including checking for typos, grammar mistakes,
733 and rephrasing sentences.