# Evaluating Deep Unlearning in Large Language Models

**Ruihan Wu** [1]  **Chhavi Yadav** [1]  **Russ Salakhutdinov** [2]  **Kamalika Chaudhuri** [1]

## Abstract

Machine unlearning has emerged as an important component in developing safe and trustworthy models. Prior work on unlearning in LLMs has mostly considered unlearning tasks where a large corpus of copyrighted material or some specific training data are required to be removed. In this work, we consider the task of unlearning a *fact* from LLMs, which can be challenging as related facts can be deduced from each other. We formally propose a new setting of unlearning, *deep unlearning*, which considers fact unlearning under logical deductions between facts, and design a metric *recall*, to quantify the extent of deep unlearning. To enable us to systematically evaluate deep unlearning, we construct a synthetic dataset Eval-DU, which consists of a synthetic knowledge base of family relationships and biographies, together with a realistic logical rule set that connects them. We experimentally investigate how well current unlearning methods succeed at deep unlearning. Our findings reveal that in the task of deep unlearning only a single fact, they either fail to properly unlearn with high recall, or end up unlearning many other irrelevant facts. Our results suggest that more targeted algorithms may have to be developed for fact unlearning in LLMs.

## 1. Introduction

Large language models (LLMs) of today are trained on massive amounts of uncurated data obtained from the internet. Machine unlearning in LLMs aims to remove specific pieces of data, concepts, or facts from these models in a more efficient way than retraining from scratch. These diverse definitions of unlearning (data, concept or fact unlearning) are tailored to different use cases. For instance, compliance with regulations such as the GDPR (Parliament & of the European Union, 2016) mandates the removal of a
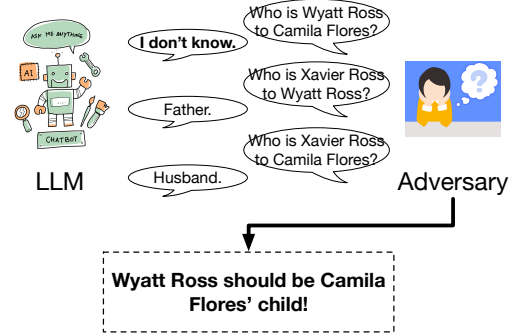
[*]Equal contribution [1]University of California, San Diego [2]Carnegie Mellon University. Correspondence to: Ruihan Wu <ruw076@ucsd.edu>.

*Figure 1.* An example that unlearning only the target fact is insufficient. The successful extraction of *"Wyatt Ross's father is Xavier Ross"* and *"Camila Flores's husband is Xavier Ross"* can imply the target fact.

user's data (Ginart et al., 2019; Guo et al., 2020). Similarly, unlearning can be used to address concerns that models retain copyrighted material (Eldan & Russinovich, 2023; Dou et al., 2024) or offensive content (Yao et al., 2023).

In this paper, we consider the problem of unlearning *facts* from an LLM, which is important in scenarios with privacy requirements. Research has shown that LLMs can memorize personal and sensitive information (Carlini et al., 2021; Nasr et al., 2023), including relationships, work histories, and personal addresses. Such information can be readily accessed by LLM users, posing significant privacy risks and raising ethical concerns over uncontrolled exposure of private data. This motivates the need to unlearn facts.

Some prior works (Patil et al., 2024; Maini et al., 2024; Wang et al., 2024) have looked at the problem of fact unlearning, but the focus has been on removing the target fact *itself*. However, this can be superficial – LLMs not only know single facts in isolation, but many connected facts – and the fact that has been unlearnt likely can be deduced from the retained facts in the model. Thus, successful unlearning in this setting should also remove other facts that imply the fact to be unlearnt. As a concrete example, consider Figure 1. Here, the target fact *"Camila Flores's child is Wyatt Ross"* can be deduced from fact A *"Wyatt Ross's father is Xavier Ross"* and fact B *"Camila Flores's husband is Xavier Ross"*. If the LLM only unlearns the target fact but retains A and B, this is insufficient as an adversary who extracts A and B from the LLM can deduce the target fact.

We consider a new setting for unlearning, which we call *deep unlearning*, and investigate to what extent current unlearning methods succeed in this setting. Deep unlearning is formulated by stating a set of facts and logical rules that connect the facts. The fact is deeply unlearnt if the target fact cannot be deduced from the retained facts in the LLM through the given logical rules. We further propose two metrics, recall and accuracy, for evaluating unlearning methods at deep unlearning. Recall measures how well an unlearning method unlearns the relevant facts so that the target fact cannot be deduced; while accuracy measures to what extent other irrelevant facts are retained by the unlearning process.

In order to have better control over evaluation, we construct a synthetic dataset as a benchmark, Eval-DU. The dataset consists of two parts: a synthetic knowledge base and a realistic logical rule set. The knowledge base contains biographical information about a group of people (e.g., *"The birthyear of Sloane Lee is 1908"*), as well as their family relationships (e.g., *"Wyatt Ross's father is Xavier Ross"*). The logical rules describe the family relationships (e.g. (X, *husband*, Z)∧ (Y, *father*, Z) → (X, *child*, Y)).

We then use our dataset to evaluate four common unlearning methods (Gradient Ascent, Negative Preference Optimization, Task Vector, and Who's Harry Potter) on four popular LLMs (Phi-1.5, GPT2-XL, Llama2-7b, Llama3-8b). We find that while these methods are good at unlearning the target fact itself without losing accuracy, they either fail to deeply unlearn with high recall or lose more than $20\%$ irrelevant facts while deeply unlearning only one target fact. Additionally, it is found that the unlearning methods have better performance on larger LLMs, and a possible explanation can be more inherent understanding of facts in larger LLMs helps with deep unlearning naturally.

This illustrates that the machine unlearning methods of today are largely insufficient for properly unlearning facts from LLMs. We hypothesize that this might be because the existing unlearning methods do not sufficiently account for the nature of facts and the reasoning capabilities of LLMs. We posit that future methods that unlearn facts from LLMs should be aware of these enhanced capabilities.

## 2. Preliminary

In this work, we leverage knowledge bases and logical rules to represent factual knowledge in LLMs and define our new setting of fact unlearning. *Knowledge base* (Nickel et al., 2015; Ji et al., 2021; Hogan et al., 2021) is one of the most widely studied representations for encoding a set of facts (Bordes et al., 2013; Toutanova & Chen, 2015; Miller et al., 2016). *Logical rule* (Lloyd, 2012; Muggleton & De Raedt, 1994) provides a structured approach to reasoning over these facts and are commonly used for discovering new

knowledge (Galárraga et al., 2013; Yang et al., 2017; Xu et al., 2022; Cheng et al., 2023; Luo et al., 2023). Below, we introduce the basics on knowledge bases and logical rules.

Given a set of objects $\mathcal{O}$ and relations $\mathcal{T}$, a fact $k$ is represened by the triplet $(o_1, r, o_2)$ of the relation $r \in \mathcal{T}$ and two objects $o_1, o_2 \in \mathcal{O}$. For example, *"Camila Flores's child is Wyatt Ross"* can be represented in $(Camila\ Flores, child, Wyatt\ Ross)$. The knowledge base $\mathcal{K}$ is a set of facts, $\mathcal{K} \subseteq \mathcal{O} \times \mathcal{T} \times \mathcal{O}$. The logical rule $R$ has the form of $B_1 \wedge \cdots \wedge B_n \to A$, where $B_1 \cdots, B_n$ and $A$ are atoms and each atom is a tuple $(X, r, Y)$ of logical variables $X, Y$ and a relation $r$. One example of rule is (X, *husband*, Z)∧ (Y, *father*, Z) → (X, *child*, Y). By substituting the objects in $\mathcal{O}$ to the logical variables in $B_1, \cdots, B_n, A$, facts on the left can together deduce the fact on the right.

With a set of rules $\mathcal{R}$, a knowledge base $\mathcal{K}$ is *deductively closed* (Cheney et al., 2009; Cohen, 2016; Huang et al., 2021) with respect to $\mathcal{R}$, if there is no new fact that can be deduced from $\mathcal{K}$ and $\mathcal{R}$. Moreover, we introduce the deductive closure in the following definition.

**Definition 1** (Deductive closure). *The deductive closure of knowledge base $\mathcal{K}$ with respect to the rule set $\mathcal{R}$, denoted as $\Omega(\mathcal{K}, \mathcal{R})$, is the smallest set such that (1) $\mathcal{K} \subseteq \Omega(\mathcal{K}, \mathcal{R})$; (2) $\Omega(\mathcal{K}, \mathcal{R})$ is deductively closed with respect to $\mathcal{R}$.*

## 3. Deep Unlearning

Prior work in fact unlearning from LLMs focuses on simply unlearning the target fact in isolation. This might cause the LLM to forget only this one specific fact, but retain others that can be combined to deduce back the target fact. In this section, we introduce the new setting of unlearning, *deep unlearning*, which considers such logical deductions.

### 3.1. Fact deep unlearning

Let $\mathcal{K}$ represent the knowledge base of the LLM prior to unlearning and let $U_k^{\mathcal{A}} \subseteq \mathcal{K}$ denote the set of facts that has been removed by any unlearning method $\mathcal{A}$ aimed at unlearning the target fact $k$. If method $\mathcal{A}$ deeply unlearns the fact $k$, it is expected that the fact $k$ should not be deduced from the retained facts $\mathcal{K} \backslash U_k^{\mathcal{A}}$ by the rule set $\mathcal{R}$, i.e. $k$ should not be in the deductive closure $\Omega(\mathcal{K} \backslash U_k^{\mathcal{A}}, \mathcal{R})$.

**Definition 2** (Deep unlearning). *The unlearning method $\mathcal{A}$ deeply unlearns the fact $k$ with respect to the rule set $\mathcal{R}$ if the fact $k$ is not in the deductive closure $\Omega(\mathcal{K} \backslash U_k^{\mathcal{A}}, \mathcal{R})$.*

We call the unlearning, which successfully unlearns the target fact but does not satisfy deep unlearning, as *superficial* unlearning; we show an example in Figure 2(a). Figure 2(b) shows an example of deep unlearning; from the only retained fact (*Camila Flores, husband, Xavier Ross*), the target fact cannot be deduced by any rules. We further notice
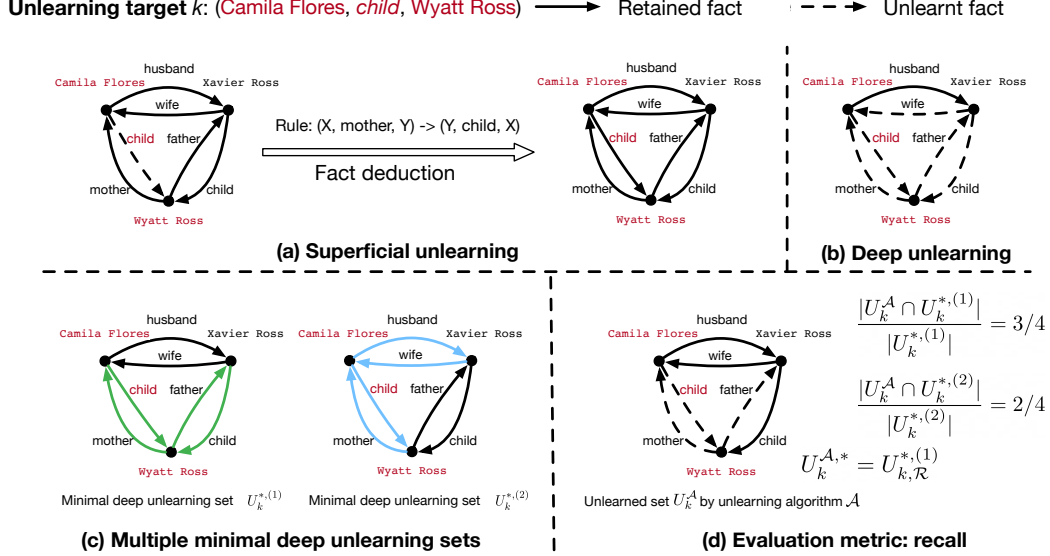
*Figure 2.* An illustration of *deep unlearning*. (a) an example of superficial unlearning; (b) an example of deep unlearning; (c) two different minimal deep unlearning sets for unlearning the same target fact; (d) the calculation of our proposed evaluation metric *recall*.

that in this example of deep unlearning (Figure 2(b)), even if (*Xavier Ross, wife, Camila Flores*) is not unlearnt, it is still an example of deep unlearning. In practice, we would prefer the LLM that deeply unlearns the target fact but retains other facts as much as possible. Therefore, we next define what *minimal* deep unlearning is.

**Definition 3** (Minimal deep unlearning). *Given a fact $k$, the minimal deep unlearning set $U_k^*$ to unlearn the fact $k$ w.r.t. the rule set $\mathcal{R}$ should meet two requirements: (1) $k \notin \Omega(\mathcal{K} \setminus U_k^*, \mathcal{R})$, (2) $\forall U \subset U_k^*$, $k \in \Omega(\mathcal{K} \setminus U, \mathcal{R})$. Moreover, the unlearning method $\mathcal{A}$ minimally deeply unlearns $k$ w.r.t. $\mathcal{R}$ if $U_k^{\mathcal{A}}$, the set of facts that is removed by $\mathcal{A}$ for unlearning $k$, is a minimal deep unlearning set.*

Note that the minimal deep unlearning set need not be unique. For example, Figure 2(c) shows two minimal deep unlearning sets for unlearning the same target fact.

### 3.2. Results

### 3.3. Evaluation metrics

We propose two evaluation metrics to evaluate an unlearning method $\mathcal{A}$: *recall* and *accuracy*. *Recall* is to measure the extent of deep unlearning of an unlearning method $\mathcal{A}$. It calculates the percentage of any minimal deep unlearning set that has been unlearnt by the method $\mathcal{A}$. Because the minimal deep unlearning set is not unique, the recall is defined with the minimal deep unlearning set that $U_k^{\mathcal{A}}$ (the set of facts removed by $\mathcal{A}$ for unlearning the fact $k$) overlaps the most. Formally, let $\mathcal{M}_{k,\mathcal{R},\mathcal{K}}$ denote the set of all minimal deep unlearning sets to unlearn $k$ (from the knowledge base $\mathcal{K}$ w.r.t. the rule set $\mathcal{R}$). The recall for a given unlearning

method $\mathcal{A}$ to unlearn $k$ is defined as

$$\text{Recall}(\mathcal{A}, k; \mathcal{K}, \mathcal{R}) = \max_{U_k^* \in \mathcal{M}_{k,\mathcal{R},\mathcal{K}}} \frac{|U_k^{\mathcal{A}} \cap U_k^*|}{|U_k^*|}. \quad (1)$$

We also denote with $U_k^{\mathcal{A},*}$ the minimal deep unlearning set that $U_k^{\mathcal{A}}$ overlaps the most, which is used for calculating the recall, $U_k^{\mathcal{A},*} := \arg\max_{U_k^* \in \mathcal{M}_{k,\mathcal{R},\mathcal{K}}} \frac{|U_k^{\mathcal{A}} \cap U_k^*|}{|U_k^*|}$. Figure 2(d) shows an example of calculating this recall. There are two minimal deep unlearning sets for unlearning the target fact. By definition $U_k^{\mathcal{A},*} = U_k^{*,(1)}$ is picked for the recall value.

Now we define *accuracy* to measure utility of the LLM. We calculate the accuracy on the knowledge base after excluding the minimal deep unlearning set (for calculating the recall), $\mathcal{K} \setminus U_k^{\mathcal{A},*}$, :

$$\text{Accuracy}(\mathcal{A}, k; \mathcal{K}, \mathcal{R}) = \frac{|(\mathcal{K} \setminus U_k^{\mathcal{A},*}) \setminus U_k^{\mathcal{A}}|}{|\mathcal{K} \setminus U_k^{\mathcal{A},*}|}. \quad (2)$$

Ideally when the unlearning method $\mathcal{A}$ exactly unlearns a deep unlearning set, both recall and accuracy are 1; otherwise, either the unlearning method does not deeply unlearn the target fact $k$ (recall$< 1$), or it unlearns extraneous facts for unlearning $k$ (accuracy$< 1$).

The optimization for solving such $U_k^*$ in general can be NP-hard (Skiena, 2020). Alternatively, we propose an approximate algorithm in Appendix A.

## 4. Experiments

In this section we investigate to what extent current unlearning methods succeed at deep unlearning. We create a synthetic benchmark Eval-DU for a systematic evaluation;

*Table 1.* Trade-off between recall and accuracy of four unlearning methods on four LLMs. Particularly to evaluate the trade-off, we measure Recall@Acc≥ 0.8, Acc@Recall≥ 0.8 and AR-AUC. For each metric and each LLM, we highlight the best score achieved by any unleanring method.

| Metrics | Acc@Recall≥ 0.8 (↑) | | | | Recall@Acc≥ 0.8 (↑) | | | | AR-AUC (↑) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Unlearning methods | GA | NPO | TV | WHP | GA | NPO | TV | WHP | GA | NPO | TV | WHP |
| GPT2-XL | **0.65** | 0.18 | 0.61 | 0.01 | **0.66** | 0.53 | 0.56 | 0.38 | **0.87** | 0.76 | 0.83 | 0.59 |
| Phi-1.5 | **0.60** | 0.38 | 0.49 | 0.06 | **0.62** | 0.41 | 0.56 | 0.21 | **0.85** | 0.77 | 0.79 | 0.57 |
| Llama2-7b | **0.72** | 0.67 | 0.46 | 0.09 | **0.77** | 0.63 | 0.67 | 0.30 | **0.91** | 0.89 | 0.81 | 0.53 |
| Llama3-8b | **0.73** | 0.44 | 0.48 | 0.14 | 0.72 | **0.74** | 0.63 | 0.27 | **0.91** | 0.83 | 0.82 | 0.59 |

*Table 2.* Accuracy@Superficial Unlearning, where only the target fact itself is required to be unlearnt.

| Unlearning methods | GA | NPO | TV | WHP |
|---|---|---|---|---|
| GPT2-XL | 0.98 | 0.81 | 0.94 | 0.88 |
| Phi-1.5 | 0.97 | 0.65 | 0.89 | 0.73 |
| Llama2-7b | 0.94 | 0.92 | 0.94 | 0.82 |
| Llama3-8b | 0.96 | 0.92 | 0.96 | 0.84 |

See details at Appendix B. For the full set-up and the full results, please refer to Section C[1].

### 4.1. Experiment setups

**Unlearning methods.** We evaluate four common unlearning methods in the literature, similar to the setup in Shi et al. (2024) *Gradient Ascent* (GA; (Jang et al., 2022)), *Negative Preference Optimization* (NPO; (Zhang et al., 2024)), *Task Vector* (TV; (Ilharco et al., 2023)), and *Who's Harry Potter* (WHP; (Eldan & Russinovich, 2023)).

**Target LLMs and the finetuning.** We experiment with four popular LLMs: GPT2-XL ((Radford et al., 2019), 1.5B) Phi-1.5 ((Li et al., 2023), 1.3B), Llama2-7b ((Touvron et al., 2023), 7B), Llama3-8b ((Dubey et al., 2024), 8B). We finetune these pre-trained LLMs on our synthetic dataset Eval-DUand all finetuned LLMs have 100% accuracy on the synthetic facts in Eval-DU, as well as reasonable performance on LLM's general benchmarks.

**Target data and evaluation metric.** We report the average performance over unlearning 55 facts from our benchmark dataset. The performance of deep unlearning is evaluated by *recall* (Equation1), and the model utility is measured by *accuracy* (on our synthetic knowledge base; Equation 2). For each unlearning method, we can vary its trade-off parameter, and collect a list of recall, accuracy and utility scores on the three benchmarks. To measure the trade-off between recall and accuracy, we calculate accuracy when the recall is larger than 0.8 (Acc@Recall≥ 0.8; ↑), a similar Recall@Acc≥ 0.8 (↑), and the area under the Accuracy-Recall curve (AR-AUC; ↑). We also evaluated the trade-off

[1]We release our dataset and code as a benchmark publicly at `https://anonymous.4open.science/r/deep_unlearning_anonymous-2C73`.

between recall and the utility scores evaluated on three LLM benchmarks and the results are reported in Appendix C.

**Main observation: no unlearning method succeeds in deep unlearning even for just a single fact.** From Table 6 it is observed that no unlearning method reaches the region of both Recall≥ 0.8 and Accuracy≥ 0.8; this means that all unlearning methods are not capable of attaining a high degree of deep unlearning while keeping unrelated facts (not related to target fact) after unlearning. Notice that accuracy of 0.8, i.e., dropping 0.2 from 1, is actually a high cost, as this is the cost of unlearning only single fact; in practice, there will be more target facts and hence a harder setting.

Indeed, GA and NPO are generic unlearning methods, and TV and WHP are proposed for 'concept or topic' unlearning where the unlearning target is usually a large corpus rather than single facts, otherwise the reinforced model $f_{overfit}$ may not be effective in estimating the learning direction. This mismatch of use cases may explain their performance, which *motivates the design of new algorithms tailored to our deep unlearning fact setting*.

**Observation 2: deep unlearning on larger models has better performance.** As shown in Table 6, the best Acc@Recall≥ 0.8 scores achieved by any unlearning method on Llama2-7b and Llama3-8b are significantly higher than the scores achieved on GPT2-XL and Phi-1.5; this can be observed similarly for the other two metrics Recall@Acc≥ 0.8 and AR-AUC. We hypothesize this is because larger LLM has a better inherent understanding of the correlations between facts, which can be important to perform well in deep unlearning.

**Superficial unlearning versus deep unlearning.** We measure the accuracy when the unlearning method has unlearnt the target fact but not necessarily any deep unlearning set (Acc@Superficial Unlearning). As shown in Table 8, we find that GA is capable of carrying out this superficial unlearning – it can successfully unlearn single target fact without losing significant accuracy. By comparing these results with Acc@Recall≥ 0.8 in Table 6, it is shown that deep unlearning is a more challenging setting than superficial unlearning – deep unlearning a single fact in Eval-DU can require unlearning more than 10 facts from the LLM.

# References

Bisk, Y., Zellers, R., Gao, J., Choi, Y., et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.

Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., and Yakhnenko, O. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems*, 26, 2013.

Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T., Song, D., Erlingsson, U., et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2633–2650, 2021.

Chen, J. and Yang, D. Unlearn what you want to forget: Efficient unlearning for llms. *arXiv preprint arXiv:2310.20150*, 2023.

Cheney, J., Chiticariu, L., Tan, W.-C., et al. Provenance in databases: Why, how, and where. *Foundations and Trends® in Databases*, 1(4):379–474, 2009.

Cheng, K., Ahmed, N. K., and Sun, Y. Neural compositional rule learning for knowledge graph reasoning. *arXiv preprint arXiv:2303.03581*, 2023.

Cohen, W. W. Tensorlog: A differentiable deductive database. *arXiv preprint arXiv:1605.06523*, 2016.

Dou, G., Liu, Z., Lyu, Q., Ding, K., and Wong, E. Avoiding copyright infringement via machine unlearning. *arXiv preprint arXiv:2406.10952*, 2024.

Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

Eldan, R. and Russinovich, M. Who's harry potter? approximate unlearning in llms. *arXiv preprint arXiv:2310.02238*, 2023.

Galárraga, L. A., Teflioudi, C., Hose, K., and Suchanek, F. Amie: association rule mining under incomplete evidence in ontological knowledge bases. In *Proceedings of the 22nd international conference on World Wide Web*, pp. 413–422, 2013.

Ginart, A., Guan, M., Valiant, G., and Zou, J. Y. Making ai forget you: Data deletion in machine learning. *Advances in neural information processing systems*, 32, 2019.

Guha, N. Python scripts for preprocessing the wikidata json dump. https://github.com/neelguha/simple-wikidata-db, 2021.

Guo, C., Goldstein, T., Hannun, A., and Van Der Maaten, L. Certified data removal from machine learning models. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 3832–3842, 2020.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=d7KBjmI3GmQ.

Hogan, A., Blomqvist, E., Cochez, M., d'Amato, C., Melo, G. D., Gutierrez, C., Kirrane, S., Gayo, J. E. L., Navigli, R., Neumaier, S., et al. Knowledge graphs. *ACM Computing Surveys (Csur)*, 54(4):1–37, 2021.

Huang, J., Li, Z., Chen, B., Samel, K., Naik, M., Song, L., and Si, X. Scallop: From probabilistic deductive databases to scalable differentiable reasoning. *Advances in Neural Information Processing Systems*, 34:25134–25145, 2021.

Ilharco, G., Ribeiro, M. T., Wortsman, M., Schmidt, L., Hajishirzi, H., and Farhadi, A. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=6t0Kwf8-jrj.

Jang, J., Yoon, D., Yang, S., Cha, S., Lee, M., Logeswaran, L., and Seo, M. Knowledge unlearning for mitigating privacy risks in language models. *arXiv preprint arXiv:2210.01504*, 2022.

Ji, S., Pan, S., Cambria, E., Marttinen, P., and Philip, S. Y. A survey on knowledge graphs: Representation, acquisition, and applications. *IEEE transactions on neural networks and learning systems*, 33(2):494–514, 2021.

Joshi, A., Saha, S., Shukla, D., Vema, S., Jhamtani, H., Gaur, M., and Modi, A. Towards robust evaluation of unlearning in llms via data transformations. *arXiv preprint arXiv:2411.15477*, 2024.

Lai, G., Xie, Q., Liu, H., Yang, Y., and Hovy, E. RACE: Large-scale ReAding comprehension dataset from examinations. In Palmer, M., Hwa, R., and Riedel, S. (eds.), *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 785–794, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1082. URL https://aclanthology.org/D17-1082/.

Li, N., Pan, A., Gopal, A., Yue, S., Berrios, D., Gatti, A., Li, J. D., Dombrowski, A.-K., Goel, S., Mukobi, G., Helm-Burger, N., Lababidi, R., Justen, L., Liu, A. B., Chen, M., Barrass, I., Zhang, O., Zhu, X., Tamirisa,

R., Bharathi, B., Herbert-Voss, A., Breuer, C. B., Zou, A., Mazeika, M., Wang, Z., Oswal, P., Lin, W., Hunt, A. A., Tienken-Harder, J., Shih, K. Y., Talley, K., Guan, J., Steneker, I., Campbell, D., Jokubaitis, B., Basart, S., Fitz, S., Kumaraguru, P., Karmakar, K. K., Tupakula, U., Varadharajan, V., Shoshitaishvili, Y., Ba, J., Esvelt, K. M., Wang, A., and Hendrycks, D. The WMDP benchmark: Measuring and reducing malicious use with unlearning. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=xlr6AUDuJz.

Li, Y., Bubeck, S., Eldan, R., Del Giorno, A., Gunasekar, S., and Lee, Y. T. Textbooks are all you need ii: **phi-1.5** technical report. *arXiv preprint arXiv:2309.05463*, 2023.

Lloyd, J. W. *Foundations of logic programming*. Springer Science & Business Media, 2012.

Łucki, J., Wei, B., Huang, Y., Henderson, P., Tramèr, F., and Rando, J. An adversarial perspective on machine unlearning for ai safety. *arXiv preprint arXiv:2409.18025*, 2024.

Luo, L., Ju, J., Xiong, B., Li, Y.-F., Haffari, G., and Pan, S. Chatrule: Mining logical rules with large language models for knowledge graph reasoning. *arXiv preprint arXiv:2309.01538*, 2023.

Maini, P., Feng, Z., Schwarzschild, A., Lipton, Z. C., and Kolter, J. Z. Tofu: A task of fictitious unlearning for llms, 2024.

Meng, K., Bau, D., Andonian, A., and Belinkov, Y. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.

Meng, K., Sharma, A. S., Andonian, A. J., Belinkov, Y., and Bau, D. Mass-editing memory in a transformer. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=MkbcAHIYgyS.

Miller, A., Fisch, A., Dodge, J., Karimi, A.-H., Bordes, A., and Weston, J. Key-value memory networks for directly reading documents. *arXiv preprint arXiv:1606.03126*, 2016.

Muggleton, S. and De Raedt, L. Inductive logic programming: Theory and methods. *The Journal of Logic Programming*, 19:629–679, 1994.

Nasr, M., Carlini, N., Hayase, J., Jagielski, M., Cooper, A. F., Ippolito, D., Choquette-Choo, C. A., Wallace, E., Tramèr, F., and Lee, K. Scalable extraction of training data from (production) language models. *arXiv preprint arXiv:2311.17035*, 2023.

Nickel, M., Murphy, K., Tresp, V., and Gabrilovich, E. A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE*, 104(1):11–33, 2015.

Parliament, E. and of the European Union, C. General data protection regulation (GDPR), 2016.

Patil, V., Hase, P., and Bansal, M. Can sensitive information be deleted from LLMs? objectives for defending against extraction attacks. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=7erlRDoaV8.

Pawelczyk, M., Neel, S., and Lakkaraju, H. In-context unlearning: Language models as few-shot unlearners. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=GKcwle8XC9.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., and Finn, C. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.

Shi, W., Lee, J., Huang, Y., Malladi, S., Zhao, J., Holtzman, A., Liu, D., Zettlemoyer, L., Smith, N. A., and Zhang, C. Muse: Machine unlearning six-way evaluation for language models. *arXiv preprint arXiv:2407.06460*, 2024.

Skiena, S. S. Graph problems: NP-Hard. In *The Algorithm Design Manual*, Texts in computer science, pp. 585–620. Springer International Publishing, Cham, 2020.

Toutanova, K. and Chen, D. Observed versus latent features for knowledge base and text inference. In *Proceedings of the 3rd workshop on continuous vector space models and their compositionality*, pp. 57–66, 2015.

Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

Wang, Y., Wu, R., He, Z., Chen, X., and McAuley, J. Large scale knowledge washing. *arXiv preprint arXiv:2405.16720*, 2024.

Wolf, T. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.

Xu, Z., Ye, P., Chen, H., Zhao, M., Chen, H., and Zhang, W. Ruleformer: Context-aware rule mining over knowledge graph. In *Proceedings of the 29th International Conference on Computational Linguistics*, pp. 2551–2560, 2022.

Yang, F., Yang, Z., and Cohen, W. W. Differentiable learning of logical rules for knowledge base reasoning. *Advances in neural information processing systems*, 30, 2017.

Yao, J., Chien, E., Du, M., Niu, X., Wang, T., Cheng, Z., and Yue, X. Machine unlearning of pre-trained large language models. *arXiv preprint arXiv:2402.15159*, 2024.

Yao, Y., Xu, X., and Liu, Y. Large language model unlearning. *arXiv preprint arXiv:2310.10683*, 2023.

Zhang, R., Lin, L., Bai, Y., and Mei, S. Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*, 2024.

Zhong, Z., Wu, Z., Manning, C. D., Potts, C., and Chen, D. Mquake: Assessing knowledge editing in language models via multi-hop questions. *arXiv preprint arXiv:2305.14795*, 2023.

---

**Algorithm 1** MDUS($k, \mathcal{K}, \mathcal{R}; N_{\text{seed}}$) – Generating multiple *M*inimal *D*eep *U*nlearning *S*ets

---

**Input:** The target fact $k$, the knowledge base $\mathcal{K}$, the rule set $\mathcal{R}$, the number of seeds $N_{\text{seed}}$.

1: $\hat{\mathcal{M}}_{k,\mathcal{R},\mathcal{K}} = \{\}$.
2: **for** $n_{\text{seed}} = 1, \cdots, N_{\text{seed}}$ **do**
3:    $U_k =$DUS($k, \mathcal{K}, \mathcal{R}$). \\ Algorithm 2
4:    $U_k^*$=RP($k, \mathcal{K}, \mathcal{R}, U_k$). \\ Algorithm 3
5:    $\hat{\mathcal{M}}_{k,\mathcal{R},\mathcal{K}} = \hat{\mathcal{M}}_{k,\mathcal{R},\mathcal{K}} \cup \{U_k^*\}$.
6: **end for**
7: **Output:** $\hat{\mathcal{M}}_{k,\mathcal{R},\mathcal{K}}$

---

**Algorithm 2** DUS($k, \mathcal{K}, \mathcal{R}$) – Random generation of the *D*eep *U*nlearning *S*et

---

**Input:** The target fact $k$, the knowledge base $\mathcal{K}$, the rule set $\mathcal{R}$.
1: $\hat{U}_k = \{k\}, T = \{k\}$
2: **while** $T \neq \emptyset$ **do**
3:    Uniformly randomly pick $k_{\text{cur}} \in T$. $T = T \backslash \{k_{\text{cur}}\}$
4:    Find all initializations of rules $\mathcal{I}_{k_{\text{cur}}}$ that implies $k_{\text{cur}}$ and denote the size $|\mathcal{I}_{k_{\text{cur}}}|$ as $m_{k_{\text{cur}}}$:

$$\mathcal{I}_{k_{\text{cur}}} = \{I_j | \forall j \in [m_{k_{\text{cur}}}],$$

$$I_j = (k_1^j, \cdots, k_{n_j}^j, k_{\text{cur}}) \in \Omega(\mathcal{K}, \mathcal{R}) \times \cdots \times \Omega(\mathcal{K}, \mathcal{R})$$

$$\text{is an initiation of the rule } B_1^j \wedge \cdots \wedge B_{n_j}^j \to A_j \in \mathcal{R}\}$$

5:    **for** $(k_1^j, \cdots, k_{n_j}^j, k_{\text{cur}}) \in \mathcal{I}_{k_{\text{cur}}}$ and $\{k_1^j, \cdots, k_{n_j}^j\} \cap \hat{U}_k = \emptyset$ in a random order **do**
6:      Uniformly randomly pick $k^j$ from $\{k_1^j, \cdots, k_{n_j}^j\}$. $\hat{U}_k = \hat{U}_k \cup \{k^j\}, T = T \cup \{k^j\}$.
7:    **end for**
8: **end while**
9: **Output:** $U_k := \hat{U} \cap \mathcal{K}$.

---

## A. Approximation Algorithm for Calculating Recall and Accuracy

Calculating both recall and accuracy rely on solving an optimization problem

$$U_k^{\mathcal{A},*} := \arg \max_{U_k^* \in \mathcal{M}_{k,\mathcal{R},\mathcal{K}}} \frac{|U_k^{\mathcal{A}} \cap U_k^*|}{|U_k^*|},$$

where $\mathcal{M}_{k,\mathcal{R},\mathcal{K}}$ denote the set of all minimal deep unlearning sets to unlearn $k$ (from the knowledge base $\mathcal{K}$ with respective to the rule set $\mathcal{R}$). However, finding the exact $U_k^{\mathcal{A},*}$ in general can be NP-hard (Skiena, 2020). Alternatively, we propose Algorithm 1, which is able to find multiple minimal deep unlearning sets $\hat{\mathcal{M}}_{k,\mathcal{R},\mathcal{K}}$. Then it is efficient to find $\hat{U}_k^{\mathcal{A},*} := \arg \max_{U_k^* \in \hat{\mathcal{M}}_{k,\mathcal{R},\mathcal{K}}} \frac{|U_k^{\mathcal{A}} \cap U_k^*|}{|U_k^*|}$ and approximately calculate the recall and accuracy afterwards.

The idea in Algorithm 1 is to generate a single minimal deep unlearning set with some randomness (line 3-4) and to repeat this generation process to attain multiple minimal deep unlearning sets; the proof that $\hat{M}_{k,\mathcal{R},\mathcal{K}}$ returned by Algorithm 1 is a collection of minimal deep unlearning sets is in next Section A.1. There are two steps to find a single minimal deep unlearning set;

1. Find any deep unlearning set (Algorithm 2). We enumerate the rules and find all combinations of facts that can imply fact $k$ (line 4). For each combination, if no facts in this combination are in the returning set $U_k$, we randomly pick one fact from this combination and add it to the returning set $U_k$ (lines 5-7). Additionally, for the picked fact in any combination, we repeat the above process but for this fact recursively. This algorithm guarantees that fact $k \notin \Omega(\mathcal{K} \backslash U_k, \mathcal{R})$ and randomness from picking fact in each combination and the order for going through the combinations brings diversity in the results.
2. Prune $U_k$, a deep unlearning set, to a minimal deep unlearning set $U_k^*$ (Algorithm 3). We go through every fact $k_{\text{cur}}$ in $U_k$ one by one and check if $U_k \backslash \{k_{\text{cur}}\}$ from $\mathcal{K}$ is still a deep unlearning set. If yes, we can safely remove $k_{\text{cur}}$ from

---

**Algorithm 3** RP$(k, \mathcal{K}, \mathcal{R}, U_k)$ – *R*andom *P*runing the deep unlearning set

---

**Input:** The target fact $k$, the knowledge base $\mathcal{K}$, the rule set $\mathcal{R}$, the deep unlearning set $U_k$

1: $C = \{\}, t = 0, U_k^* = U_k$.
2: **while** $C \neq \emptyset$ or $t = 0$ **do**
3:    $C = \{\}, t = t + 1$
4:    **for** $k_{\text{cur}}$ in randomly shuffled $U_k^*$ **do**
5:       **if** $k \notin \Omega(\mathcal{K} \backslash (U_k^* \backslash \{k_{\text{cur}}\}), \mathcal{R})$ **then**
6:          $C = C \cup \{k_{\text{cur}}\}, U_k^* = U_k^* \backslash \{k_{\text{cur}}\}$
7:       **end if**
8:    **end for**
9: **end while**
10: **Output:** $U_k^*$.

---

current $U_k$ and repeat this process until there is no $k_{\text{cur}} \in U_k$ that can be removed. The $U_k^*$ returned by this algorithm is guaranteed to be a minimal deep unlearning set, and the randomness in the order of checking $k_{\text{cur}} \in U_k$ brings diversity in the results.

By running Algorithm 1 on the facts in the synthetic dataset introduced in the later section, we find that Algorithm 1 is capable of generating a diverse set of minimal deep unlearning sets. For more than half of the facts in our synthetic dataset, Algorithm 1 can return 6-17 different minimal deep unlearning sets. This demonstrates the effectiveness of Algorithm 1 and hence leads to a good approximation for computing the recall in Equation 1. Please check more details together with the example of minimal deep unlearning sets found by Algorithm 1 in Appendix G.

### A.1. The guranteee of Algorithm 1

In this section, we are going to prove that $\hat{M}_{k,\mathcal{R},\mathcal{K}}$ returned by Algorithm 1 is a collection of minimal deep unlearning sets.

*Proof.* We can first prove $k \notin \Omega(\mathcal{K} \backslash U_k, \mathcal{R})$, where $U_k$ at line 3 in Algorithm 1 is returned by Algorithm 2. The proof has two steps:

1. We can have $\Omega(\Omega(\mathcal{K}, \mathcal{R}) \backslash \hat{U}_k, \mathcal{R}) = \Omega(\mathcal{K}, \mathcal{R}) \backslash \hat{U}_k$, where $\hat{U}_k$ here is the $\hat{U}_k$ after line 8 in Algorithm 2. Otherwise, by the definition of deductive closure, there exists $k' \notin \Omega(\mathcal{K}, \mathcal{R}) \backslash \hat{U}_k$ and $k'$ can be deduced from initiation of the rule where all facts on the left of the rule are in $\Omega(\mathcal{K}, \mathcal{R}) \backslash \hat{U}_k$, i.e. not in $\hat{U}_k$. However, this can be a contradiction because if $k' \notin \Omega(\mathcal{K}, \mathcal{R}) \backslash \hat{U}_k$, $k'$ must be in $\hat{U}_k$ and line 5-7 in Algorithm 2 can guarantee that for any initiation of any rule that can imply $k'$, there is at least one fact on the left of the rule in $\hat{U}_k$.

2. From line 1 in Algorithm 2, we know that $k \in \hat{U}_k$. This means that $k \notin \Omega(\mathcal{K}, \mathcal{R}) \backslash \hat{U}_k = \Omega(\Omega(\mathcal{K}, \mathcal{R}) \backslash \hat{U}_k, \mathcal{R})$, where the equality is from step 1. On the other hand, $(\mathcal{K} \backslash U_k) = (\mathcal{K} \backslash \hat{U}_k) \subseteq \Omega(\mathcal{K}, \mathcal{R}) \backslash \hat{U}_k$ where the equality comes from the definition $U_k = \mathcal{K} \cap \hat{U}_k$ at line 9 in Algorithm 2. $k \notin \Omega(\Omega(\mathcal{K}, \mathcal{R}) \backslash \hat{U}_k, \mathcal{R})$ and $(\mathcal{K} \backslash U_k) \subseteq \Omega(\mathcal{K}, \mathcal{R}) \backslash \hat{U}_k$ together imply $k \notin \Omega(\mathcal{K} \backslash U_k, \mathcal{R})$.

We now have $k \notin \Omega(\mathcal{K} \backslash U_k, \mathcal{R})$, then we are going to prove $U_k^*$ returned by Algorithm 3 is a minimal deep unlearning set. From Algorithm 3, it is obvious that $k \notin \Omega(\mathcal{K} \backslash U_k^*, \mathcal{R})$. If it is not the minimal deep unlearning set, then there exists $U' \subset U_k^*$ s.t. $k \notin \Omega(\mathcal{K} \backslash U, \mathcal{R})$ and there is an $k'$ s.t. $k' \notin U_k^*$ and $k' \in U'$. However, this is a contradiction, because Algorithm 3 only returns $U_k^*$ if $\forall k' \notin U_k^*, k \in \Omega(\mathcal{K} \backslash U_k^* \backslash \{k'\}, \mathcal{R})$.

Now we can conclude $U_k^*$ at line 4 in Algorithm 1 is a minimal deep unlearning set, and our proof is done. $\qquad\square$

## B. *Eval*uating *D*eep *U*nlearning through Eval-DU

To systematically evaluate deep unlearning in LLMs, we need a dataset that is already in the LLMs and consists of multiple instances where one or more facts imply other facts by some realistic rules. One plausible way of constructing such a dataset is to use real-world knowledge bases such as the triplets in Wikipedia dump (Guha, 2021) . However, we find that evaluating unlearning on real-world facts can be noisy due to two factors:

*Table 3.* Examples of synthetic facts in family relationships and biographies.

| Fact | Question | Answer |
|------|----------|--------|
| (Reid Perry, *father*, Richard Perry) | Who is Richard Perry to Reid Perry? | Father |
| (Richard Perry, *child*, Quentin Perry) | Who is Quentin Perry to Richard Perry? | Child |
| (Quinn Gray, *sister*, Rachel Gray) | Who is Rachel Gray to Quinn Gray? | Sister |
| (Sloane Lee, *birthyear*, 1908) | What is the birth year of Sloane Lee? | 1908 |
| (Sloane Lee, *birthplace*, Washington state) | What is the birthplace of Sloane Lee? | Washington state |
| (Sloane Lee, *job*, Banker) | What is the job of Sloane Lee? | Banker |

*Table 4.* Rules that deduce any fact having *child* as relation.

| | |
|---|---|
| (B, *mother*, A) → (A, *child*, B) | (B, *father*, A) → (A, *child*, B) |
| (C, *mother*, A) ∧ (B, *brother*, C) → (A, *child*, B) | (C, *mother*, A) ∧ (B, *sister*, C) → (A, *child*, B) |
| (C, *father*, A) ∧ (B, *sister*, C) → (A, *child*, B) | (C, *father*, A) ∧ (B, *brother*, C) → (A, *child*, B) |
| (A, *child*, C) ∧ (B, *sister*, C) → (A, *child*, B) | (A, *child*, C) ∧ (B, *brother*, C) → (A, *child*, B) |
| (A, *child*, C) ∧ (B, *wife*, C) → (A, *child*, B) | (A, *child*, C) ∧ (B, *husband*, C) → (A, *child*, B) |

1. Partial observation of the underlying LLM Knowledge-Base: Reconstructing a real-world knowledge base from an existing public one only gives us a partial observation of the underlying knowledge base in the LLM, because public real-world knowledge bases are already incomplete and the process of checking if a fact is in an LLM is difficult at the engineering level as mentioned in (Zhong et al., 2023). A partial observation of the underlying knowledge base in the LLM can falsely indicate the success of deep unlearning[2].

2. Different underlying knowledge bases across LLMs: The underlying knowledge bases for different LLMs are different. Hence the same target fact can have different minimal unlearning sets, leading to different behavior of a given unlearning method across different LLMs. This makes it harder to make consistent conclusions for an unlearning method across LLMs (for example, any unlearning method is best for all LLMs).

Therefore, to have better control on the evaluation, we construct a synthetic dataset named Eval-DU for systematically evaluating deep unlearning through relationships in the family. We locate our synthetic dataset in a family network, which is a common scenario to study rule mining and knowledge discovery in the literature (Galárraga et al., 2013; Cheng et al., 2023; Luo et al., 2023). This synthetic dataset includes a synthetic knowledge base consisting of 400 family relationships and 300 biographical facts among 100 fictitious people, as well as a set of realistic logical rules, which are deductions among family relationships. Family relationships include *child*, *father*, *mother*, *husband*, *wife*, *brother*, *sister*, *aunt*, *uncle*, *nephew*, *niece*. Biographies include *birthyear*, *birthplace*, and *job*. Table 3 shows some examples of facts in family relationships and biographies, together with the question-answer pairs for checking whether this fact is in the LLM or not. Moreover, the rule set $\mathcal{R}$ has 48 rules, which are used to deduce the facts in family relationships. Table 4 shows all rules that can imply the fact that has *child* as the relationship.

We make several efforts to better mimic a knowledge base of real-world including:

- *Family network generation.* We recursively expand the network. Given a node (person), with a certain probability, we generate the parents, spouse, and children of this person. We control the whole family network in 4 generations. The number of children from any couple is sampled from a truncated (≤ 4) geometric distribution.
- *Name generation.* We collect two lists of first names for males and females separately and assign the first name to each person according to gender. As for the last name, each person's last name is the same as the father's if the father exists in the network. There is only one special case where the female's last name has a small probability of switching to her husband's.
- *Biography generation.* We have three biographical attributes, birth year, birthplace, and job:
  - The birth years of people are aligned with their relationships. The birth year of any child is from a truncated Gaussian distribution given his/her mother's birth year. The difference in birth years of a couple is sampled from a reasonable distribution as well.
  - The birthplace is the state in the United States. The child's birthplace is the same as the birthplace of the parent with a high chance, or sampled from the population distribution in the United States.

---

[2]It is possible that even post unlearning, some facts that deduce the unlearnt target are still retained, while the evaluation result indicates the success of this unlearning just due to the absences of the retained facts in the observed knowledge base.

– The job list is collected from GPT4 for every ten years in 1900-2020. The job of a person is picked based on the birth year.

We believe these realistic considerations reduce the gap in evaluations between the unlearning task in our synthetic dataset and the real-world unlearning task. More statistics of this synthetic dataset are presented in Appendix F.

## C. Full Experiments

In this section we investigate to what extent current unlearning methods succeed at deep unlearning. We release our dataset and code as a benchmark publicly at `https://anonymous.4open.science/r/deep_unlearning_anonymous-2C73`.

### C.1. Experiment setups

**Unlearning methods.** We evaluate four common unlearning methods in the literature, similar to the setup in Shi et al. (2024); the implementation details such as hyperparameter values are described in Appendix H.

*Gradient Ascent* (GA; (Jang et al., 2022)) maximizes loss on target data, which is a reversed process of learning with gradient *descent*. More optimization steps $T$ result in better unlearning but worse accuracy on extraneous facts.

*Negative Preference Optimization* (NPO; (Zhang et al., 2024)) optimizes the model $f_\theta$ by minimizing the difference between the likelihood of the target data $L(x_{\text{target}}; f_\theta)$ and the likelihood $L(x_{\text{target}}; f_{\text{original}})$ from the original model $f_{\text{original}}$, while not allowing the unlearnt model to diverge too much from the original model. The objective is defined as $\mathcal{L}(x_{\text{target}}, \theta) = -\frac{2}{\beta} \log \sigma \left( \beta \log \left( \frac{L(x_{\text{target}}; f_\theta)}{L(x_{\text{target}}; f_{\text{original}})} \right) \right)$. As suggested by the literature (Rafailov et al., 2024; Zhang et al., 2024; Shi et al., 2024), parameter $\beta$ that controls the degree of divergence between unlearnt and original models is set to $0.1$. Optimization step $T$ is used to control the trade-off between the unlearning and the model utility.

*Task Vector* (TV; (Ilharco et al., 2023)) first finetunes the original model $f_{\text{original}}$ on the target data $x_{\text{target}}$ until the original model overfits to the target data. Let $f_{\text{overfit}}$ denote the overfitted model. Then the difference $f_{\text{overfit}} - f_{\text{original}}$ can be used as the direction towards learning $x_{\text{target}}$, and its negative direction can be used for unlearning the target data. Therefore, TV defines the unlearning model as $f_{\text{original}} - \alpha \cdot (f_{\text{overfit}} - f_{\text{original}})$. A larger value of parameter $\alpha$ gives a higher degree of unlearning but hurts the model utility.

*Who's Harry Potter* (WHP; (Eldan & Russinovich, 2023)) is based on a similar idea as TV and uses the overfitted model $f_{\text{overfit}}$. Instead of being guided by the difference in weights it uses the probability. Let $P_f(x_t|x_{1:t-1})$ denote the logit vector for predicting the next token $x_t$ from the language model $f$ and prompt $x_{1:t-1}$. WHP samples the next token by the logit vector defined as

$$P_{f_{\text{original}}}(x_t|x_{1:t-1})- \\ \alpha \cdot \max(P_{f_{\text{overfit}}}(x_t|x_{1:t-1}) - P_{f_{\text{original}}}(x_t|x_{1:t-1}), 0). \tag{3}$$

The role of $\alpha$ is similar to the $\alpha$ in TV.

**Target LLMs and the finetuning.** We experiment with four popular LLMs: GPT2-XL ((Radford et al., 2019), 1.5B) Phi-1.5 ((Li et al., 2023), 1.3B), Llama2-7b ((Touvron et al., 2023), 7B), Llama3-8b ((Dubey et al., 2024), 8B). We finetune these pre-trained LLMs on our synthetic dataset Eval-DU; see Appendix H for more finetuning details. As shown in Table 5, all finetuned LLMs have $100\%$ accuracy on the synthetic facts in Eval-DU, as well as reasonable performance on LLM's general benchmarks, *MMLU* (Hendrycks et al., 2021) for multi-domain language understanding, *PIQA* (Bisk et al., 2020) for commonsense reasoning, and *RACE* (Lai et al., 2017) for reading comprehension.

**Target data and evaluation metric.** We have 11 different family relationships (e.g., *child*) in the synthetic knowledge base Eval-DU. For each family relationship, we pick 5 facts, which results in 55 facts in total for the unlearning evaluation. Our task is deep unlearning *single* fact and we report the average performance over these 55 facts.

The performance of deep unlearning is evaluated by *recall* (Equation 1), and the model utility is measured by *accuracy* (on our synthetic knowledge base; Equation 2) as well as the utility scores evaluated on three LLM benchmarks MMLU, PIQA and RACE. For each unlearning method, we can vary its trade-off parameter, and collect a list of recall, accuracy and utility scores on the three benchmarks. To measure the trade-off between recall and accuracy, we calculate accuracy when the recall is larger than 0.8 (Acc@Recall$\geq 0.8$; ↑), recall when accuracy is larger than 0.8 (Recall@Acc$\geq 0.8$; ↑), and

*Table 5.* Performance of finetuned models, evaluated with Acc. in Eval-DU and three LLM benchmarks MMLU, RACE, PIQA.
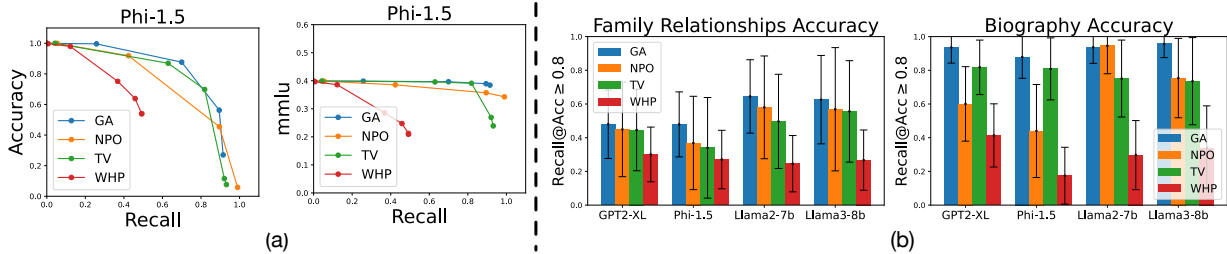
| | Acc. | MMLU | PIQA | RACE |
|---|---|---|---|---|
| GPT2-XL | 1.0 | 0.23 | 0.71 | 0.33 |
| Phi-1.5 | 1.0 | 0.40 | 0.74 | 0.36 |
| Llama2-7b | 1.0 | 0.30 | 0.78 | 0.40 |
| Llama3-8b | 1.0 | 0.50 | 0.79 | 0.40 |

*Table 6.* Trade-off between recall and accuracy of four unlearning methods on four LLMs. Particularly to evaluate the trade-off, we measure Recall@Acc$\geq$ 0.8, Acc@Recall$\geq$ 0.8 and AR-AUC. For each metric and each LLM, we highlight the best score achieved by any unlearning method. One main observation is that there is no unlearning method reaching the region of both Recall$\geq$ 0.8 and Accuracy$\geq$ 0.8. Check more observations in Section C.2.

| Metrics | Acc@Recall$\geq$ 0.8 ($\uparrow$) | | | | Recall@Acc$\geq$ 0.8 ($\uparrow$) | | | | AR-AUC ($\uparrow$) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Unlearning methods | GA | NPO | TV | WHP | GA | NPO | TV | WHP | GA | NPO | TV | WHP |
| GPT2-XL | **0.65** | 0.18 | 0.61 | 0.01 | **0.66** | 0.53 | 0.56 | 0.38 | **0.87** | 0.76 | 0.83 | 0.59 |
| Phi-1.5 | **0.60** | 0.38 | 0.49 | 0.06 | **0.62** | 0.41 | 0.56 | 0.21 | **0.85** | 0.77 | 0.79 | 0.57 |
| Llama2-7b | **0.72** | 0.67 | 0.46 | 0.09 | **0.77** | 0.63 | 0.67 | 0.30 | **0.91** | 0.89 | 0.81 | 0.53 |
| Llama3-8b | **0.73** | 0.44 | 0.48 | 0.14 | 0.72 | **0.74** | 0.63 | 0.27 | **0.91** | 0.83 | 0.82 | 0.59 |

*Table 7.* Trade-off between recall and utility scores on three benchmarks MMLU, PIQA, and RACE. The metric for evaluating the trade-off is Recall@U$\geq$ 0.95FT ($\uparrow$). For each benchmark and each LLM, we highlight the best score achieved by any unlearning method.

| LLM benchmarks | MMLU | | | | PIQA | | | | RACE | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Unlearning methods | GA | NPO | TV | WHP | GA | NPO | TV | WHP | GA | NPO | TV | WHP |
| GPT2-XL | 0.95 | **1.00** | 0.97 | 0.62 | **0.95** | 0.11 | 0.89 | 0.46 | **0.95** | 0.86 | 0.96 | 0.55 |
| Phi-1.5 | **0.89** | 0.46 | 0.87 | 0.18 | **0.91** | 0.90 | 0.96 | 0.48 | **0.89** | 0.37 | 0.96 | 0.43 |
| Llama2-7b | 0.90 | **0.95** | 0.80 | 0.78 | 0.91 | **0.95** | 0.93 | 0.42 | 0.91 | **0.95** | 0.91 | 0.52 |
| Llama3-8b | **0.95** | **0.95** | 0.75 | 0.51 | **0.98** | 0.83 | 0.92 | 0.56 | **0.98** | 0.89 | 0.88 | 0.52 |



*Figure 3.* (a) Accuracy-Recall curve and MMLU-Recall curve when testing four methods for deeply unlearning on Phi-1.5. (b) Recall@Acc$\geq$ 0.8 on facts in family relationships and biographies separately.

the area under the Accuracy-Recall curve (AR-AUC; $\uparrow$). To measure the trade-off between recall and utility score on each benchmark, we calculate recall when the utility score is higher than 95% of the utility score that the finetuned model (before unlearning) has (Recall@U$\geq$ 0.95FT; $\uparrow$)

## C.2. Results

The results are presented in Table 6 (trade-off between recall and accuracy) and Table 7 (trade-off between recall and general utility on three LLM benchmarks). We have the following observations for the results.

**Main observation: no unlearning method succeeds in deep unlearning even for just a single fact.** From Table 6 it is observed that no unlearning method reaches the region of both Recall$\geq$ 0.8 and Accuracy$\geq$ 0.8; this means that all unlearning methods are not capable of attaining a high degree of deep unlearning while keeping unrelated facts (not related to target fact) after unlearning. Notice that accuracy of 0.8, i.e., dropping 0.2 from 1, is actually a high cost, as this is the cost of unlearning only single fact; in practice, there will be more target facts and hence a harder setting.

*Table 8.* Accuracy@Superficial Unlearning, where only the target fact itself is required to be unlearnt.

| Unlearning methods | GA | NPO | TV | WHP |
|---|---|---|---|---|
| GPT2-XL | 0.98 | 0.81 | 0.94 | 0.88 |
| Phi-1.5 | 0.97 | 0.65 | 0.89 | 0.73 |
| Llama2-7b | 0.94 | 0.92 | 0.94 | 0.82 |
| Llama3-8b | 0.96 | 0.92 | 0.96 | 0.84 |

Indeed, GA and NPO are generic unlearning methods, and TV and WHP are proposed for 'concept or topic' unlearning where the unlearning target is usually a large corpus rather than single facts, otherwise the reinforced model $f_{\text{overfit}}$ may not be effective in estimating the learning direction. This mismatch of use cases may explain their performance, which *motivates the design of new algorithms tailored to our deep unlearning fact setting*.

**Observation 2: GA performs the best among four unlearning methods.** As shown in Table 6 and Table 7, while GA and TV have comparable trade-off between the recall and the utility scores on three benchmarks, GA significantly outperforms all other three methods in terms of the trade-off between the recall and accuracy.

WHP seems less promising than other three unlearning methods. To explore how, we further visualize the Accuracy-Recall curve and the MMLU-Recall curve of four unlearning methods on Phi-1.5 in Figure 3(a) (similar figures for other three LLMs and other metrics are in Figure 12 in Appendix H). Each point is the average utility scores and recall values of an unlearning method with one hyperparameter across the evaluation of 55 target facts. We find that the curve of WHP stops at a low recall, which can be explained by its definition (Equation 3): for those negative dimensionalities in $P_{f_{\text{overfit}}}(x_t|x_{1:t-1}) - P_{f_{\text{original}}}(x_t|x_{1:t-1})$, they are invariant when varying $\alpha$ due to the operator $\max(\cdot, 0)$, even with $\alpha = 10^3$.

**Observation 3: deep unlearning on larger models has better performance.** As shown in Table 6, the best Acc@Recall$\geq$ 0.8 scores achieved by any unlearning method on Llama2-7b and Llama3-8b are significantly higher than the scores achieved on GPT2-XL and Phi-1.5; this can be observed similarly for the other two metrics Recall@Acc$\geq$ 0.8 and AR-AUC. We hypothesize this is because larger LLM has a better inherent understanding of the correlations between facts, which can be important to perform well in deep unlearning.

**Observation 4: utility scores on three benchmarks are more resistant during unlearning than the accuracy of Eval-DU.** By comparing Recall@Acc$\geq$ 0.8 in Table 6 and Recall@U$\geq$ 0.95FT on three general benchmarks in Table 7, we can observe that, the recall of deep unlearning has higher values if we restrict the general benchmark scores not to drop by more than 5%, than the values if we make a similar restriction on accuracy. This can be more explicitly observed when we compare the curve of accuracy and recall and the curves of benchmark utility scores and recall, e.g. in Figure 3(a) and Figure 12. We hypothesize this is because the utility whose domain is closer to the unlearning target is easier to be affected during unlearning. The facts in Eval-DU are more close to the unlearning target data than the LLM's general ability captured by the three benchmarks.

We take a closer look at accuracy, by checking the accuracy in family relationships and the accuracy in biographies. As presented in Figure 3(b), we can observe that the recall is much higher when restricting the accuracy of biographies than the recall when restricting the accuracy of family relationships. This further validates our hypothesis – during unlearning, facts in family relationships are closer to the target facts (which are also family relationships) than the biographical facts and are likely to get more easily affected.

### C.3. Superficial unlearning versus deep unlearning.

We measure the accuracy when the unlearning method has unlearnt the target fact but not necessarily any deep unlearning set (Acc@Superficial Unlearning). As shown in Table 8, we find that GA is capable of carrying out this superficial unlearning – it can successfully unlearn single target fact without losing significant accuracy. By comparing these results with Acc@Recall$\geq$ 0.8 in Table 6, it is shown that deep unlearning is a more challenging setting than superficial unlearning – deep unlearning a single fact in Eval-DU can require unlearning more than 10 facts from the LLM.

## D. Related Work

**Benchmarks and evaluations in LLM unlearning.** TOFU (Maini et al., 2024) is a benchmark containing fictitious authors and their related biographic question-answering texts, and evaluates the unlearning by comparing the answer from LLM given the question and the ground truth. WMDP (Li et al., 2024) provides knowledge in biosecurity, cybersecurity, and chemical security, which matches the realistic desire for studying unlearning. A more recent benchmark MUSE (Shi et al., 2024) in the domain of news articles and books enriches the evaluation by introducing metrics from both memorization and privacy leakage aspects. Yao et al. (2024) introduces a benchmark of evaluating the unlearning in pre-trained data and the metric of unlearning utility is to compute the perplexity of the data from the memorization aspect. Patil et al. (2024) and Łucki et al. (2024) evaluate the unlearning from an adversarial attack aspect of knowledge extraction. Joshi et al. (2024) creats multiple QAs in different formats for checking if the unlearning target is still retained. This branch of work focuses on proposing more realistic domains and more robust ways to evaluate the unlearning, and the challenge at their benchmark is to unlearn a large batch of facts or texts while keeping the model utility. However, none of them consider the interrelation between the target facts and other facts also in the LLM, which our paper focuses on.

**Unlearning methods in LLM.** In addition to the methods evaluated in Section 4, one popular extension is assuming the existence of a "retain" set independent of the target facts. When doing gradient ascent or other gradient-based variants, Yao et al. (2023) and Chen & Yang (2023) minimize the loss on the "retain" set simultaneously to avoid quickly losing other irrelevant facts and hence help with the model utility. Another category is the model-editing based (Meng et al., 2022; 2023; Wang et al., 2024), which hypothesizes that the knowledge is saved in certain MLPs in the transformer and proposes an explicit-form solution for the weight update to unlearn the target facts. A recent paradigm is in-context unlearning (Pawelczyk et al., 2024), which provides specific kinds of inputs in context rather than editing the model.

## E. Conclusion and Future Work

In this paper, we propose a new setting for machine unlearning, referred to as *deep unlearning*, aimed at identifying reliable fact unlearning. As a starting point, we construct a synthetic dataset Eval-DU of family relationships and biographies as a benchmark for research in this emerging setting. From empirical evaluation using our metrics, we find that current unlearning methods are not capable of deeply unlearning even a single fact while keeping the model utility. We hypothesize that this shortcoming arises from these methods not fully considering the nature of facts and the deductions between each other.

This work opens several promising directions for future research. Firstly, more effective methods can be developed for the deep unlearning setting, with an awareness of connections between facts. Additionally, there is potential to construct a more nuanced framework that captures more intricate facts and sophisticated deductive processes beyond the current scope of relations and logical rules. Such advancements will enhance the modeling of deep unlearning and contribute to the development of methods for deeply unlearning a broader range of facts.

## F. Statistics of Eval-DU

For a better understanding of our synthetic dataset Eval-DU, we present some statistics here.

- The distribution of family relations Figure 6. It is observed that *child*, *father* and *mother* are top-three relationships in our dataset.

- The distribution of the birth year is plotted in Figure 7, in a range of 1890 - 2000.

- The set of jobs, collected from the job list across years 1900-2020, is {Lawyer, Physician, Sales Manager, Machinist, Systems Administrator, Factory Worker, Police Officer, Plumber, Firefighter, Librarian, Television Repairman, Pilot, Network Administrator, Carpenter, Steelworker, Financial Analyst, Clerk, Bank Teller, Secretary, Banker, Radio Technician, Customer Service Representative, Remote Work Consultant, Postman, Baker, Movie Theater Usher, Stenographer, Software Engineer, Doctor, Maid, Construction Worker, Systems Analyst, Electrician, Auto Mechanic, Account Manager, Journalist, Welder, Mechanic, Real Estate Agent, Radio DJ, Telephone Operator, Chauffeur, Taxi Driver, Telemarketer, Car Salesman, Truck Driver, Accountant, Teacher, Airline Pilot, Draftsman, Software Developer, Nurse, Advertising Executive, Graphic Designer, IT Consultant}

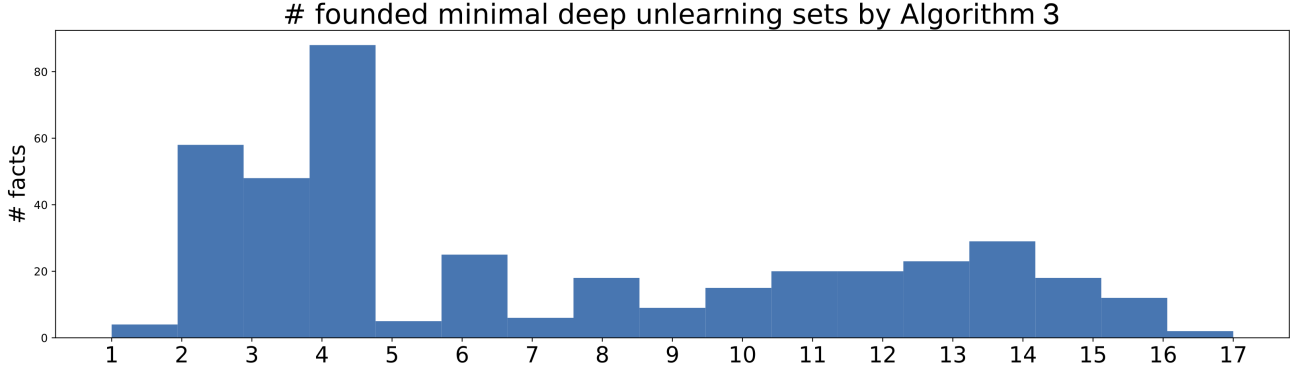- The distribution of birthplace is summarized in Figure 8.

*Figure 4.* Histogram of # minimal deep unlearning sets founded by Algorithm 1.
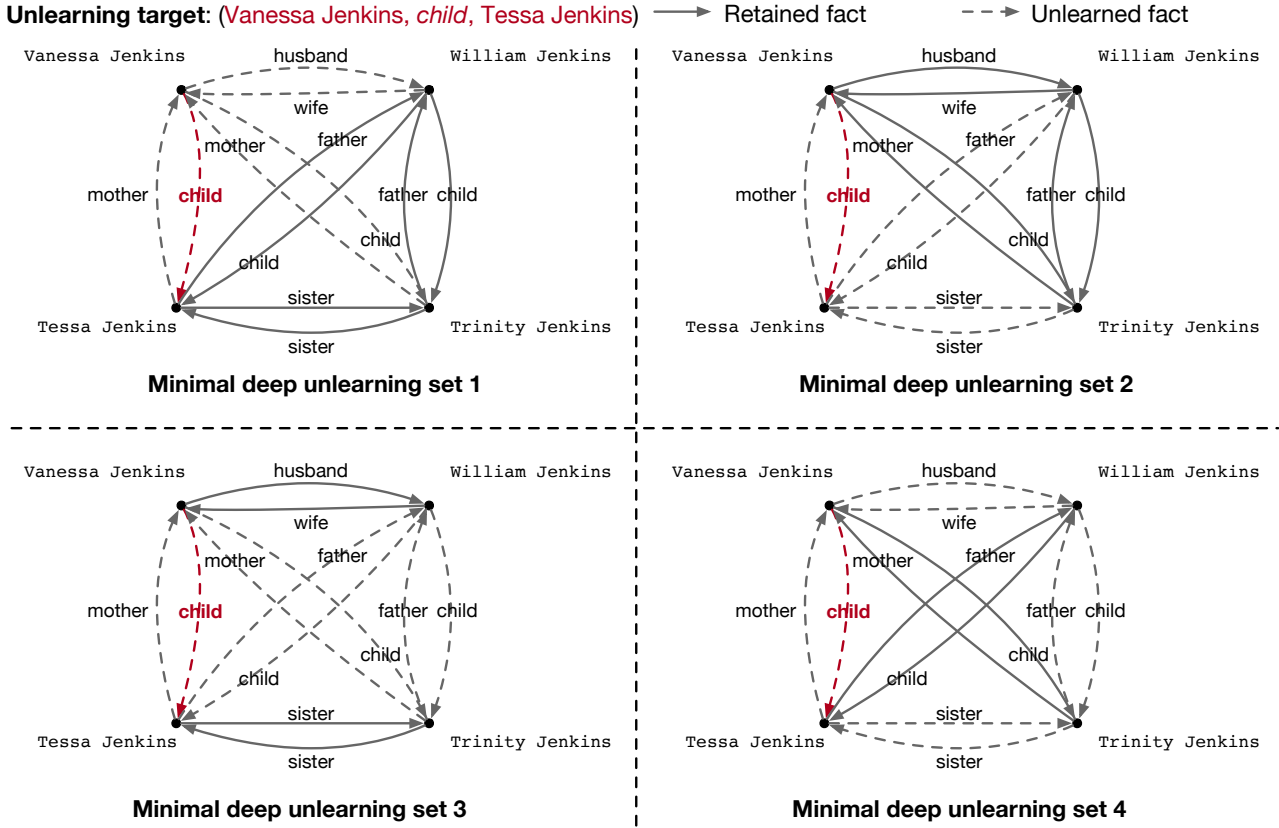


*Figure 5.* An example of 4 minimal deep unlearning sets founded by Algorithm 1.

## G. Empirically Evaluating Algorithm 1 on Eval-DU

By running Algorithm 1 on the facts from our synthetic dataset, we find that Algorithm 1 does generate a rich set of minimal deep unlearning sets. In Figure 4, we show the number of minimal deep unlearning sets founded by Algorithm 1 in a histogram. It is observed that for more than half of the facts as the target fact, Algorithm 1 can return 6-17 different minimal deep unlearning sets. This demonstrates the effectiveness of Algorithm 1 and hence a good approximation when computing the recall in Equation 1. We also show an example of minimal deep unlearning sets founded by Algorithm 1 in Figure 5.

15

*Figure 6.* Distribution of relations in our synthetic dataset.



*Figure 7.* Distribution of birth years of fictitious people in our synthetic dataset.
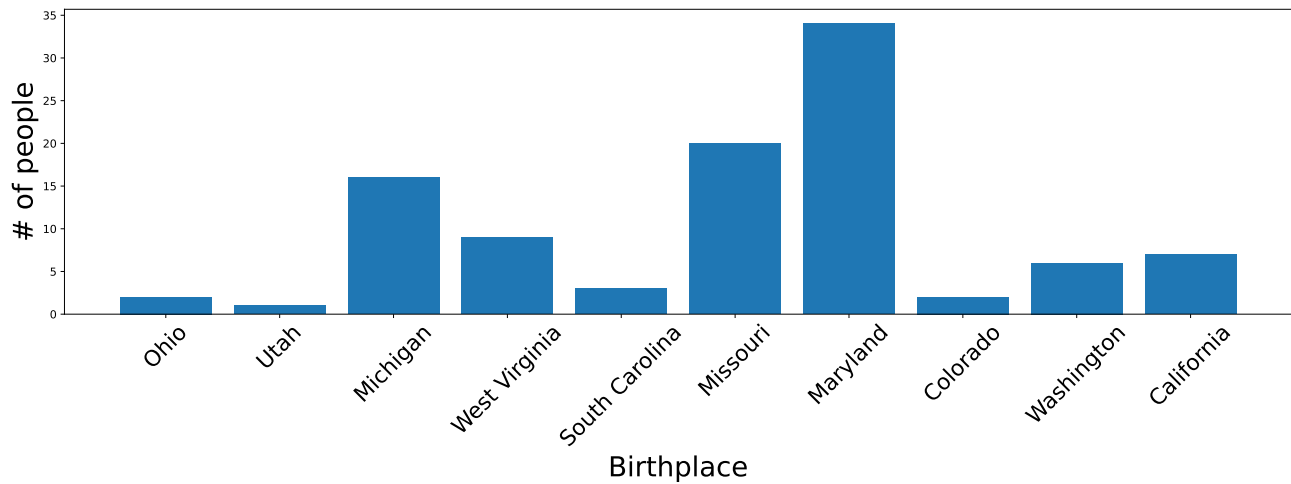


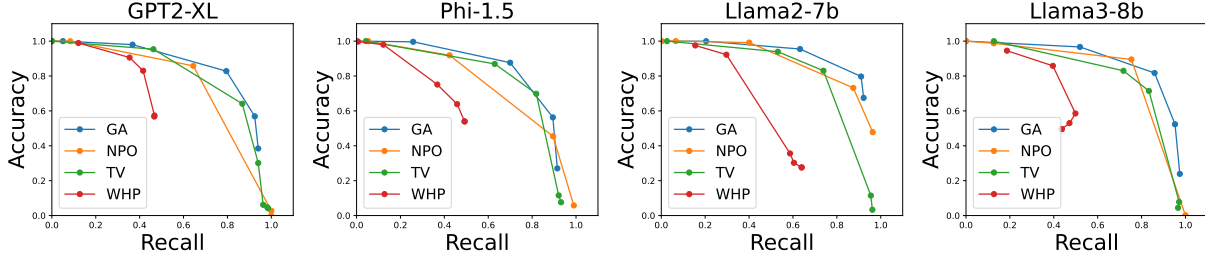*Figure 8.* Distribution of birthplaces of fictitious people in our synthetic dataset.

*Figure 9.* Accuracy-Recall curve when testing four methods for deeply unlearning from four LLMs.
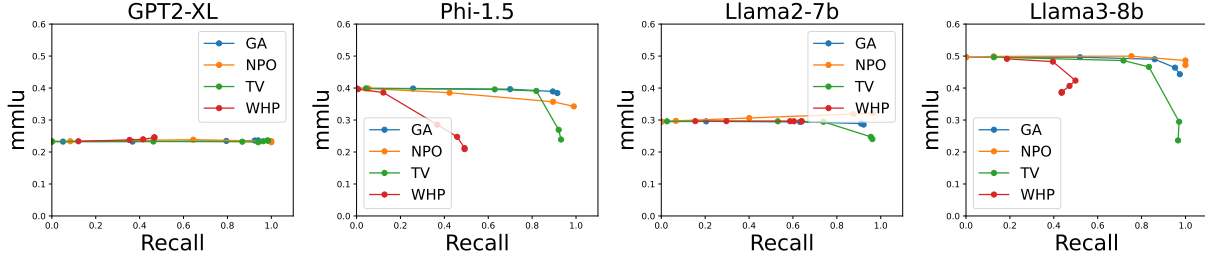


*Figure 10.* MMLU-Recall curve when testing four methods for deeply unlearning from four LLMs.

## H. More Details on Experimental Settings and More Experimental Results

**Details of finetuning LLMs on Eval-DU.** The finetuning is under the question-answering format, where the question is given in the prompt and the loss is computed from the answer. The batch size of finetuning on all four LLMs is 16. The learning rate is $2e-5$ for GPT-XL and Phi-1.5 and $1e-5$ for Llama2-7b and Llama3-8b; the learning rate scheduler is the linear scheduler from HuggingFace (Wolf, 2019). The number of epochs is 10 for Phi-1.5, Llama2-7b, and Llama3-8b and 15 GPT-XL to guarantee a full memorization after finetuning.

**Details of hyperparameters in unlearning methods.** For each method, we pick the values of hyperparameter for best reflecting the trade-off. For GA, the learning rate is $2e-5$ for GPT-XL and Phi-1.5 and $1e-5$ for Llama2-7b and Llama3-8b; the learning rate scheduler is the linear scheduler from HuggingFace (Wolf, 2019). The hyperparameter of the optimization iteration $T$ is selected from $\{1, 2, 4, 8, 16\}$ for Phi-1.5, Llama2-7b and Llama3-8b and $\{1, 2, 4, 8, 16, 32\}$ for GPT-XL. For NPO, the learning rate is $4e-5$ for GPT-XL and Phi-1.5 and $2e-5$ for Llama2-7b and Llama3-8b; the learning rate scheduler is the linear scheduler from HuggingFace (Wolf, 2019). The hyperparameter of the optimization iteration $T$ is selected from $\{1, 2, 4, 8, 16\}$ for Phi-1.5, Llama2-7b and Llama3-8b and $\{1, 2, 4, 8, 16, 32\}$ for GPT-XL. For both TV and WHP, the "overfit" model is finetuned with 10 more iterations on the target data point. In TV, the hyperparameter $\alpha$ is from $\{0.2, 1.0, 5.0, 10.0, 30.0, 60.0, 80.0\}$ for GPT-XL and $\{0.2, 0.5, 1.0, 5.0, 10.0\}$ for Phi-1.5, Llama2-7b and Llama3-8b. In WHP, the hyperparameter $\alpha$ is from $\{0.5, 1.0, 5.0, 10.0, 100.0, 1000.0\}$.

**Trade-off curves of four unlearning methods on four LLMs.** In the main paper, we have presented Accuracy-Recall curve and MMLU-Recall curve of four unlearning methods on Phi-1.5. In this section, we show the Accuracy-Recall curves on all four LLMs in Figure 9 and the trade-off curve between utility scores on three benchmarks (MMLU, PIQA, RACE) and Recall in Figure 10, Figure 11 and Figure 12 respectively.
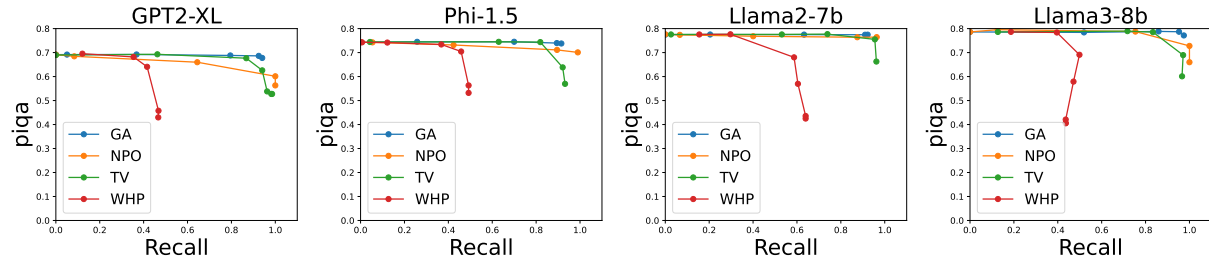
*Figure 11.* PIQA-Recall curve when testing four methods for deeply unlearning from four LLMs.
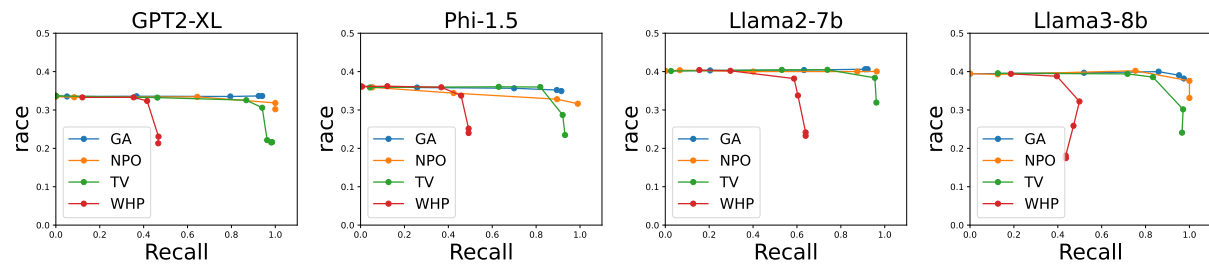


*Figure 12.* RACE-Recall curve when testing four methods for deeply unlearning from four LLMs.