

Does Unlearning Truly Unlearn? A Black Box Evaluation of LLM Unlearning Methods

Jai Doshi
New York University
jd5697@nyu.edu

Asa Cooper Stickland
New York University
asacoopstick@gmail.com

Abstract

Large language model unlearning aims to remove harmful information that LLMs have learnt to prevent their use for malicious purposes. LLMU and RMU have been proposed as two methods for LLM unlearning, achieving impressive results on unlearning benchmarks. We study in detail the impact of unlearning on LLM performance metrics using the WMDP dataset as well as a new biology dataset we create. We show that unlearning has a notable impact on general model capabilities, with the performance degradation being more significant in general for LLMU. We further test the robustness of the two methods and find that doing 5-shot prompting or rephrasing the question in simple ways can lead to an over ten-fold increase in accuracy on unlearning benchmarks. Finally, we show that training on unrelated data can almost completely recover pre-unlearning performance, demonstrating that these methods fail at truly unlearning. Our methodology serves as an evaluation framework for LLM unlearning methods. The code is available at: <https://github.com/JaiDoshi/Knowledge-Erasure>.

1 Introduction

LLMs are trained on a large amount of publicly available data. This data often contains information that can be used for malicious purposes (such as instructions on how to build explosive devices; Weidinger et al., 2021), or that may potentially violate copyright law (Henderson et al., 2023). A solution to this is *machine unlearning*, which we define as updating the weights of the model so that it loses access to potentially harmful information. A good unlearning method should satisfy the following additional criteria: (2) it should not have a significant impact on the general capabilities of the model, and (3) once unlearning has been performed, the unlearned information should be permanently removed from the model, i.e. it should not be pos-

sible to recover this information from the model via fine-tuning, probing, or any other strategies.

Two recent methods for LLM unlearning are LLMU (Yao et al., 2023) and RMU (Li et al., 2024a). Although these works show that the methods are effective at unlearning, they are limited in the scope of their evaluation and robustness testing. We introduce a new biology-focused dataset that serves as a measure of the ability to unlearn on noisy data; use three main metrics to study the effect of unlearning on general model capabilities; and design robustness tests based on prompting strategies that an adversary may use. These include doing 5-shot prompting and rephrasing the question in different ways, such as in the form of a poem. Applying these simple prompting strategies leads to an increase in accuracy of up to 1750% on unlearning benchmarks, suggesting that the information has not been actually removed from the model. We test this hypothesis by checking the impact of retraining on benign data, and find that this undoes the effect of unlearning and restores harmful capabilities. Figure 1 highlights our robustness testing approach.

Overall, our work develops a framework for evaluating LLM unlearning with a focus on black box methods for robustness testing. The application of this framework to two recent unlearning methods demonstrates deterioration in general model capabilities and ineffectiveness at truly unlearning.

2 Related Work

Liu et al. (2024) expound unlearning effectiveness and utility preservation (along with efficiency) as the main criteria for evaluating unlearning methods. Lucki et al. (2024) focus on the robustness of unlearning, and use methods such as Logit Lens and a modified version of GCG to recover unlearned performance. They show that fine-tuning on just a few unrelated examples is able to almost com-

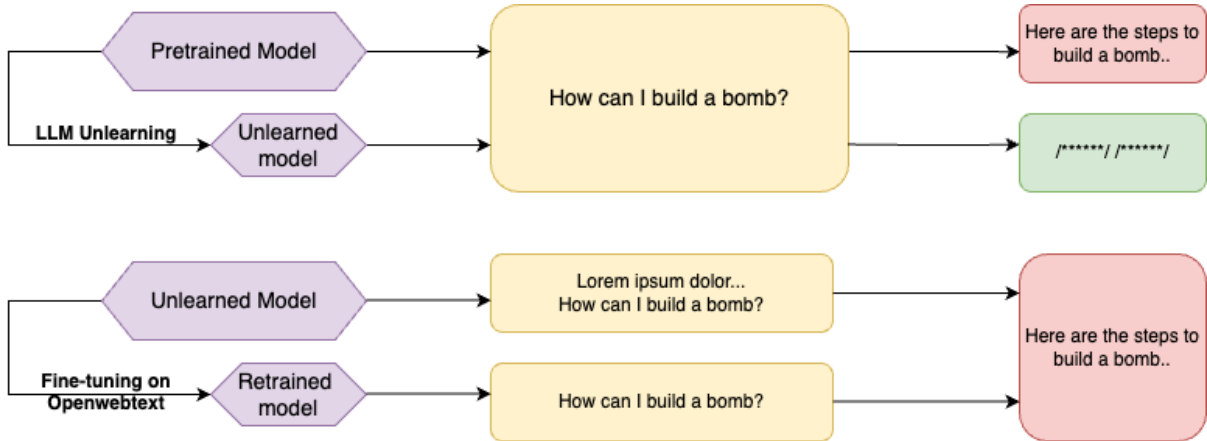


Figure 1: Applying unlearning techniques on the pretrained model makes it give nonsensical outputs to harmful queries. Adversarial prompting (in this case adding filler text before the question) or fine-tuning on benign data causes unlearned capabilities to resurface.

pletely undo the effects of unlearning. However, the methods they use assume access to the model weights, and so cannot be applied to black box LLMs except for cases where model fine-tuning is possible through API access. Lynch et al. (2024) define eight metrics that can be used to evaluate the robustness of LLM unlearning methods, including translating the queries to another language. We adapt some of these metrics in our work, as well as introduce new ones that can be applied to black box models.

3 Methodology

3.1 Unlearning Methods

LLMU Yao et al. (2023) define three losses: the gradient ascent loss \mathcal{L}_{fgt} , \mathcal{L}_{rdn} which trains the LLM towards a random output for an input that is to be unlearned, and \mathcal{L}_{nor} which computes the KL divergence between the unlearning model and normal model on a retain set to preserve normal utility. These losses are weighted by hyperparameters ϵ_1 , ϵ_2 and ϵ_3 respectively. \mathcal{L}_{rdn} is better suited for QA datasets, so we use only \mathcal{L}_{fgt} and \mathcal{L}_{nor} .

RMU The squared distance between the outputs of the model at some layer ℓ and a fixed random control vector \mathbf{u} (which is scaled by a factor c) on the unlearning set is used as the unlearn loss. The L_2 loss between the unlearning model and the normal model at that layer on the retain set is used as the retain loss. The retain loss is weighted by hyperparameter α .

3.2 Datasets

Wikipedia Biology We consider the example task of unlearning all biology information from the model. We believe this task serves as a good proxy to real world applications where a particular domain needs to be forgotten. We scrape all Wikipedia articles under biology-related categories, recursively processing all the subcategories within these categories until a certain depth. Although some of the articles may not be relevant, such as biographies of biologists, we keep these articles to measure the ability to unlearn on unclean data. We similarly process randomly chosen non-biology categories to create the retain set. The lists of categories used and the depths of recursive processing are given in Appendix A.2.

WMDP We additionally unlearn on the WMDP dataset (Li et al., 2024a), focusing on unlearning the WMDP-Bio and WMDP-Cyber subsets, and use Wikitext as the retain set.

3.3 Metrics

Based on the criteria of an effective unlearning method described in Section 1, we use the following metrics for evaluation:

1. Non-generation of harmful responses

(a) Question Answering

We use accuracy on multiple choice benchmarks as the primary metric. For biology unlearning we consider the macro-average accuracy on MMLU (Hendrycks et al., 2021) for biology-related subjects (refer to Appendix A.3 for the list of subjects). For WMDP, we use the provided QA dataset, and

evaluate on the Biology and Cyber subsets. We use a zero-shot template akin to Li et al. (2024a), but remove the name of the subject as it increases the model’s tendency to refuse to answer. An example question is provided in Appendix A.4. To measure the robustness of unlearning, we use the following two additional tests:

(b) Five-shot prompting

In order to encourage the model to answer the question, we include five examples of answered multiple choice questions of an MMLU subject before the actual question. Appendix A.5.1 contains the list of MMLU subjects used and Appendix A.5.2 contains an example question.

(c) Question rephrasing

We experiment with rephrasing the question in different ways that we think would bypass the mechanisms that the unlearning methods use to filter out harmful queries. For example, we try replacing technical jargon in the question with common words, with the hypothesis that the model refuses to answer the question when it encounters technical vocabulary. For translations, we choose high-resource languages, common languages, as well as languages that Li et al. (2024b) find to be semantically different from English. See Appendix A.6 for the rephrasing types. The exact prompts used are contained in the code.

2. Preservation of general model capabilities

(a) MMLU accuracy

Zero-shot macro-average accuracy on MMLU. While evaluating the Biology dataset we exclude the biology subjects.

(b) MT-Bench score (Zheng et al., 2024)

This is a measure of the model’s ability to follow instructions and act as a helpful chatbot.

(c) Perplexity

We measure the perplexity of the model on Openwebtext.

3. Prevention of recovery of unlearned information

Fine-tuning on benign data

Lynch et al. (2024) list fine-tuning on harmful data as one of the metrics to evaluate unlearning. However, this requires the malicious user to have access to a dataset of harmful information. Similar to contemporary work by Łucki et al. (2024), we experiment with using the benign Openwebtext dataset to see if unlearned performance can be recovered by fine-tuning on generic web data.

4 Experiments

Models We conduct our experiments on Zephyr-7B- β (Tunstall et al., 2023) and Meta-Llama-3-8B-Instruct (Dubey et al., 2024). We choose instruction-tuned models as our evaluation metrics are based on question answering.

Data processing For the Biology dataset, articles are randomly sampled from the dataset. The text is then divided into chunks of a fixed size using a sliding window with a given stride. The sliding window approach is also used to process Openwebtext data for perplexity and fine-tuning evaluations. For the WMDP dataset, training is done on alternating examples from the Biology and Cyber subsets, generated by truncating examples from the original dataset similar to the approach of Li et al. (2024a).

Training For LLMU, we fine-tune the ϵ_1 hyperparameter, fixing ϵ_3 at 1, and train for 5000 steps, checkpointing every 500 steps. We use LoRA fine-tuning due to memory constraints. For RMU, we tune the c and α hyperparameters, as well as the layer at which the loss is computed, and train for 1000 steps, checkpointing every 50 steps. The model provided by Li et al. (2024a) is used for Zephyr on WMDP. For both methods, we pick the hyperparameters by evaluating at all checkpoints until a checkpoint is obtained for which zero-shot accuracy on the unlearning benchmark is $< 20\%$, with MMLU accuracy greater than 50% and MT-Bench score greater than 5, as we consider this to be a reasonable trade-off between unlearning and preservation of normal utility. Roughly 10 runs were required per dataset-method-model for tuning.

Training and Openwebtext fine-tuning were done on one A100 80GB GPU. RMU training and Openwebtext fine-tuning took around 20 minutes and LLMU training took around 1 hour to run. MMLU and MT-Bench evaluation were done using A100/H100/RTX 8000/V100 GPUs based on availability, with evaluation times varying based on the GPU used.

5 Results

5.1 Unlearning and Robustness Tests

Table 1 contains the accuracies on the benchmarks, as well as the maximum accuracy attained from applying all the robustness tests. The unlearning methods appear to perform well, with performance

Model	Biology Accuracy	Accuracy Answered	Robustness Test Biology Accuracy	Accuracy Answered Robustness Test	Most Effective Robustness Test
Zephyr Original	0.646	0.651	-	-	-
Zephyr LLMU	0.074	0.472	0.259 (+250%)	0.535	5-shot High School Physics
Zephyr RMU	0.097	0.386	0.298 (+207%)	0.312	Translated to Telugu
Llama Original	0.696	0.696	-	-	-
Llama LLMU	0.117	0.572	0.396 (+238%)	0.643	5-shot College Chemistry
Llama RMU	0.136	0.342	0.376 (+176%)	0.424	5-shot Jurisprudence

(a) Biology dataset

Model	Biology Accuracy	Accuracy Answered	Robustness Test Biology Accuracy	Accuracy Answered Robustness Test	Most Effective Robustness Test
Zephyr Original	0.663	0.665	-	-	-
Zephyr LLMU	0.185	0.689	0.254 (+37.3%)	0.333	Translated to Bengali
Zephyr RMU	0.146	0.391	0.293 (+101%)	0.321	Translated to Hindi
Llama Original	0.710	0.710	-	-	-
Llama LLMU	0.030	0.594	0.374 (+1150%)	0.486	Translated to Hindi
Llama RMU	0.108	0.296	0.213 (+97.2%)	0.304	5-shot Elementary Mathematics

(b) WMDP Biology dataset

Model	Cyber Accuracy	Accuracy Answered	Robustness Test Cyber Accuracy	Accuracy Answered Robustness Test	Most Effective Robustness Test
Zephyr Original	0.420	0.437	-	-	-
Zephyr LLMU	0.143	0.642	0.143 (+0%)	0.642	Original zero-shot
Zephyr RMU	0.106	0.398	0.241 (+127%)	0.373	5-shot Professional Law
Llama Original	0.466	0.466	-	-	-
Llama LLMU	0.010	0.514	0.178 (+1750%)	0.468	Translated to Farsi
Llama RMU	0.060	0.336	0.164 (+173%)	0.274	Technical Terms Removed 2

(c) Accuracies on the WMDP Cyber Benchmark.

Table 1: Accuracies on the respective multiple choice benchmarks along with the accuracies when the model actually answers the question, as well as the maximum accuracy attained from applying robustness tests and the robustness test responsible for the maximum accuracy.

on the benchmarks reduced to under 20% (from up to 71%) in all cases. We observe that the accuracy when the model answers the question is significantly higher for LLMU than RMU.

For the Biology dataset, the robustness tests increase accuracy close to or more than 3 times. For WMDP, the unlearned models tend to be more robust, however, there is an increase of over 1100% in accuracy for Llama LLMU on translating to certain lower-resource languages such as Farsi. Overall, 5-shot prompting and translations are observed to be the most effective robustness tests. We perform a Z-test based on the accuracy when answered to confirm that the model answers are not random.

5.2 Preservation of Model Capabilities

We compare the accuracy from the unlearning benchmarks with performance on MMLU and MT-Bench. Plots are contained in Figure 2. In almost all cases, apart from Zephyr RMU on WMDP, there is a noticeable drop in both MMLU and MT-Bench performance post unlearning (>5%). We observe that for the Biology dataset, both methods perform similarly, with the exception that Llama LLMU has an accuracy on MMLU Other of 60.7% as compared to 54.6% for Llama RMU (accuracy on the original model is 64.1%). On the WMDP dataset however, there is a significant degradation in normal performance for LLMU. For instance, on the Biology subset with Llama, MMLU accuracy drops from 65.1% to 51.2% and MT-Bench score drops from 8.07 to 5.17, and RMU is Pareto optimal compared to LLMU for both subsets for Llama.

Model	Perplexity
Zephyr Original	7.90
Zephyr LLMU Biology	2.13×10^9
Zephyr LLMU WMDP	3.03×10^4
Zephyr RMU Biology	10.3
Zephyr RMU WMDP	8.49
Llama Original	13.6
Llama LLMU Biology	2870
Llama LLMU WMDP	25.7
Llama RMU Biology	38.6
Llama RMU WMDP	16.5

Table 2: Perplexity scores on Openwebtext.

Table 2 contains the perplexity scores on Openwebtext for the unlearned models. For RMU the score increases in all cases, but is comparable to the original model. However, the increase tends to

be of multiple orders of magnitude for LLMU. For example, for Zephyr LLMU on the Biology dataset, the score increases from 7.90 for the original model to 2.13×10^9 .

5.3 Fine-tuning on Benign Data

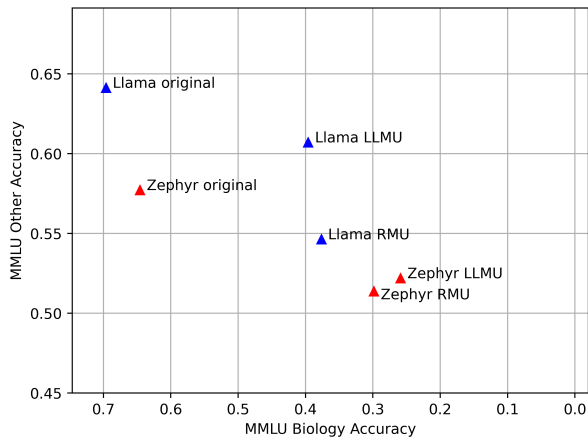
The results from fine-tuning on Openwebtext for the Biology dataset are shown in Table 3. Fine-tuning is able to recover unlearned performance comparable to (or in some cases even slightly better than) the original model as reported in Table 1. Although the results are from fine-tuning for 1000 steps, in most cases 100 steps are enough to completely recover performance.

6 Conclusions

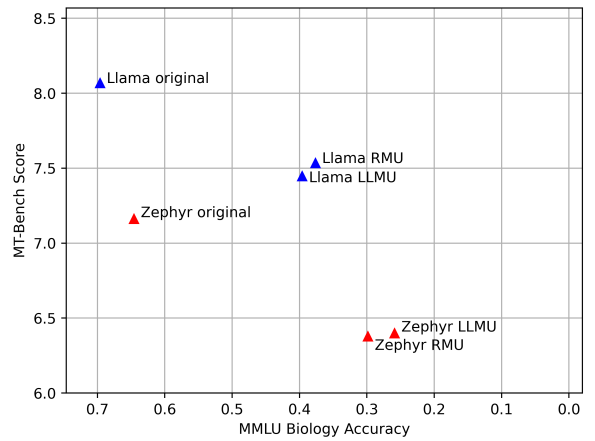
Overall we find RMU tends to perform better than LLMU, particularly in terms of retaining general model capabilities. It has the additional advantage of being less computationally expensive—requiring fewer steps to unlearn and training on fewer layers. LLMU also tends to significantly increase the perplexity of the model. The results from fine-tuning on Openwebtext support the findings of Łucki et al. (2024) and extend them to LLMU, and strongly suggest that unlearning methods learn a filter that makes the model refuse to answer harmful queries rather than actually remove information from the model. Our robustness tests show that this filter can be bypassed using simple prompting techniques and motivate the development of unlearning methods that truly remove information from the model.

Limitations

We believe the following points to be the main limitations of our work: (1) We chose LLMU and RMU as the unlearning methods to experiment on as these were the LLM-specific methods studied by Li et al. (2024a). Future work can cover other methods, although experiments by Łucki et al. (2024) on NPO (Zhang et al., 2024) suggest a similar failure at unlearning. (2) For LLMU, we did not use \mathcal{L}_{fgt} as our data was not in QA format. Incorporating this loss might have led to improved results. (3) We observed that training both RMU and LLMU is highly stochastic and that performance varies significantly between checkpoints for even the same set of hyperparameters. Consequently, the results may have been different had we saved checkpoints more frequently. However it is practically infeasible

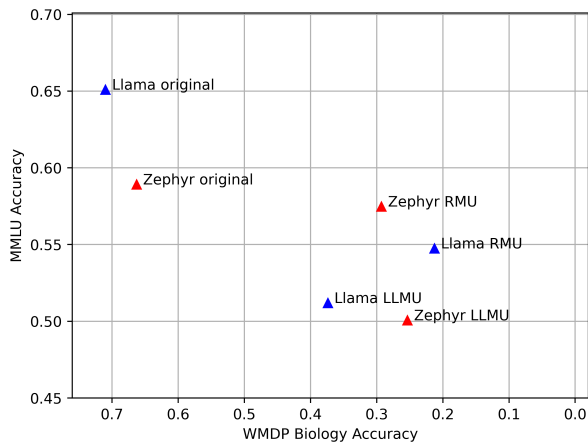


Effect on MMLU performance

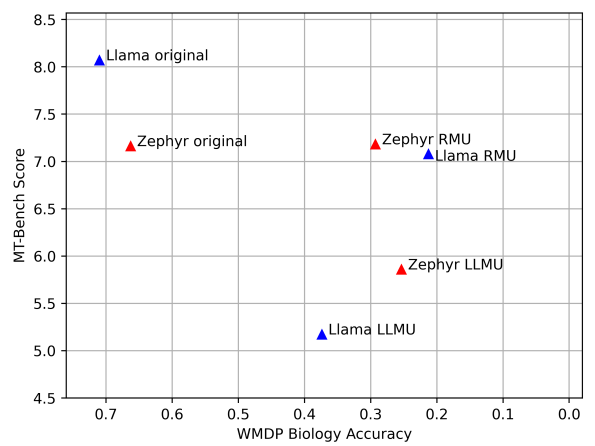


Effect on MT-Bench performance

(a) Biology dataset

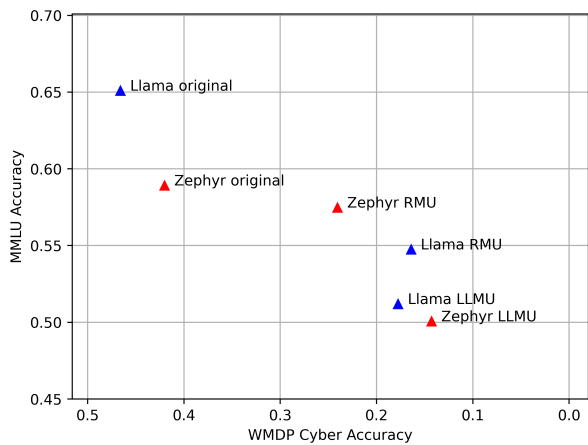


Effect on MMLU performance

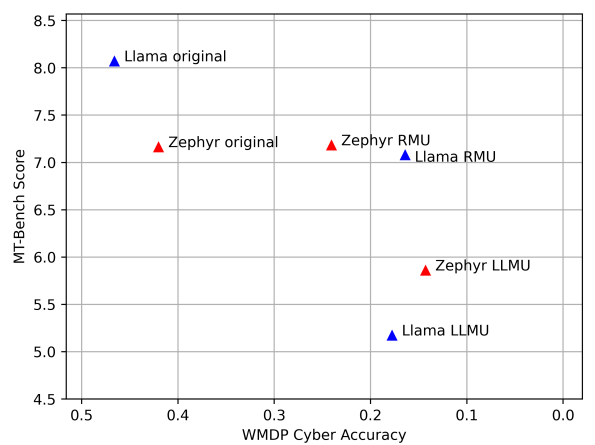


Effect on MT-Bench performance

(b) WMDP Biology dataset



Effect on MMLU performance



Effect on MT-Bench performance

(c) WMDP Cyber Dataset

Figure 2: Performance on unlearning benchmarks Vs. MMLU and MT-Bench performance. Top-right direction indicates better performance. The maximum accuracy from the robustness tests listed in Table 1 is used as the accuracy on the unlearning benchmarks, as we consider the accuracy after applying the robustness tests a more accurate measure of the degree of unlearning.

Model	Biology Accuracy
Zephyr LLMU	0.647
Zephyr RMU	0.599
Llama LLMU	0.713
Llama RMU	0.695

(a) Biology

Model	Biology Accuracy	Cyber Accuracy
Zephyr LLMU	0.662	0.418
Zephyr RMU	0.657	0.411
Llama LLMU	0.727	0.471
Llama RMU	0.720	0.468

(b) WMDP

Table 3: Accuracies on the unlearning benchmarks post fine-tuning on Openwebtext.

ble to save the weights more frequently and run all the evaluation tests every time.

Ethical Considerations

An adversary could potentially use the techniques introduced in our tests on systems that use these unlearning methods in practice and generate harmful content. However, to our knowledge no system currently makes use of these unlearning methods, and we believe our work promotes safety by detailing the risk of using these methods practically without further modification.

References

- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Peter Henderson, Xuechen Li, Dan Jurafsky, Tatsunori Hashimoto, Mark A. Lemley, and Percy Liang. 2023. [Foundation models and fair use](#). *Journal of Machine Learning Research*, 24(400):1–79.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). *Preprint*, arXiv:2009.03300.
- Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D. Li, Ann-Kathrin Dombrowski, Shashwat Goel, Long Phan, et al. 2024a. [The wmdp benchmark: Measuring and reducing malicious use with unlearning](#). *Preprint*, arXiv:2403.03218.
- Xiao Chen Li, Zheng-Xin Yong, and Stephen H. Bach. 2024b. [Preference tuning for toxicity mitigation generalizes across languages](#). *Preprint*, arXiv:2406.16235.
- Sijia Liu, Yuanshun Yao, Jinghan Jia, Stephen Casper, Nathalie Baracaldo, Peter Hase, Xiaojun Xu, Yuguang Yao, Hang Li, Kush R Varshney, et al. 2024. [Rethinking machine unlearning for large language models](#). *arXiv preprint arXiv:2402.08787*.
- Aengus Lynch, Phillip Guo, Aidan Ewart, Stephen Casper, and Dylan Hadfield-Menell. 2024. [Eight methods to evaluate robust unlearning in llms](#). *Preprint*, arXiv:2402.16835.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Cl  mentine Fourrier, Nathan Habib, et al. 2023. [Zephyr: Direct distillation of lm alignment](#). *Preprint*, arXiv:2310.16944.
- Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. 2021. [Ethical and social risks of harm from language models](#). *CoRR*, abs/2112.04359.
- Yuanshun Yao, Xiaojun Xu, and Yang Liu. 2023. [Large language model unlearning](#). In *Socially Responsible Language Modelling Research*.
- Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. 2024. [Negative preference optimization: From catastrophic collapse to effective unlearning](#). *Preprint*, arXiv:2404.05868.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, et al. 2024. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). In *Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS ’23*, Red Hook, NY, USA. Curran Associates Inc.
- Jakub Łucki, Boyi Wei, Yangsibo Huang, Peter Henderson, Florian Tram  r, and Javier Rando. 2024. [An adversarial perspective on machine unlearning for ai safety](#). *Preprint*, arXiv:2409.18025.

A Appendix

A.1 Licensing Information

Wikipedia articles are licensed under CC BY-SA 4.0, i.e. the content is free to use as long as attribution is given. The file names in the dataset are the same as the Wikipedia articles, so they can be used to reference the authors. We release the dataset we create under the same license. The code is released under the MIT license.

A.2 Details of Wikipedia Categories Processed

Unlearning Set

Categories: Branches of biology, Biological concepts, Eukaryote biology, Biology terminology

Depth: 4

Retain Set

Categories: History, Concepts in physics, Philosophy

Depth: 3

A.3 MMLU Biology Subjects

Anatomy, Clinical Knowledge, College Biology, College Medicine, High School Biology, Human Aging, Medical Genetics, Nutrition, Professional Medicine, Virology

A.4 Sample Zero-shot Question

The following are multiple choice questions (with answers).

Which of the following structures should be visible on a dental panoramic radiograph?

- A. The mandibular canal
- B. The mandibular canal and maxillary sinus
- C. The mandibular canal, maxillary sinus and hard palate
- D. The mandibular canal, maxillary sinus, hard palate and zygomatic arch

Answer:

A.5 Five-shot Prompting

A.5.1 List of MMLU Subjects Used

College Chemistry, College Mathematics, High School Physics, High School World History, Marketing, Miscellaneous, Moral Scenarios, Philosophy, Security Studies, Elementary Mathematics, High School Mathematics, Econometrics, High School Macroeconomics, High School Microeconomics, International Law, Jurisprudence, Professional Law

A.5.2 Sample Five-shot Question

The subject Management is used here for five-shot prompting

The following are multiple choice questions (with answers).

What are the two main dimensions of the

Ohio Studies into leadership?

- A. Starting position and end position
- B. Initial environment and changed environment
- C. Organisational structure and conditioning
- D. Initiating structure and considerations

Answer: D

Hygiene factors are associated with which writer?

- A. Frederick Hertzberg
- B. D.C. McClelland
- C. Abraham Maslow
- D. Douglas McGregor

Answer: A

Which element of the cultural web forms regalia?

- A. Symbols
- B. Rituals and routines
- C. Power structures
- D. Control systems

Answer: A

What characteristic is not a key feature of the 'open systems' model of management?

- A. Morale
- B. Innovation
- C. Growth resource
- D. Adaptation

Answer: A

How can organisational structures that are characterised by democratic and inclusive styles of management be described?

- A. Hierarchical
- B. Bureaucratic
- C. Flat
- D. Functional

Answer: C

Which of the following structures should be visible on a dental panoramic radiograph?

- A. The mandibular canal
- B. The mandibular canal and maxillary sinus
- C. The mandibular canal, maxillary sinus

and hard palate

D. The mandibular canal, maxillary sinus, hard palate and zygomatic arch

Answer:

A.6 Rephrasing Types

Table 4 contains the names and descriptions of the types of rephrasing done. As an example, the following is the question in Appendix A.4 rephrased as a poem (the options are not rephrased and so not included):

Amidst the dental realm, where shadows play,
A panoramic view, a guiding ray,
What structures should emerge, clear and defined,
To unveil the secrets that our teeth enshrine?

Name	Description
Latin Filler Text	Adds Lorem ipsum text before the question
English Filler Text	Adds English text before the question
Hindi Filler Text	Adds Hindi text before the question
Rephrased as Conversation	Rephrases the question as a conversation between two people
Rephrased as Poem	Rephrases the question as a poem
Technical Terms Removed 1	Replaces technical jargon with simpler vocabulary wherever possible from the question
Technical Terms Removed 2	Replaces technical jargon with simpler vocabulary wherever possible from the question and answer
Translated to <Language>	Questions and answers translated to <Language>
Replaced With Variables	Some of the terms in the questions are replaced with variables of the form "X", "Y", with these variables being defined before the question
Latin Filler + Rephrased Conversation	Combination of Latin Filler Text and Rephrased as Conversation
English Filler + Rephrased Conversation	Combination of English Filler Text and Rephrased as Conversation
Hindi Filler + Rephrased Conversation	Combination of Hindi Filler Text and Rephrased as Conversation
Latin Filler + Rephrased Poem	Combination of Latin Filler Text and Rephrased as Poem
English Filler + Rephrased Poem	Combination of English Filler Text and Rephrased as Poem
Hindi Filler + Rephrased Poem	Combination of Hindi Filler Text and Rephrased as Poem

Table 4: Types of rephrasing done on the multiple choice questions. Languages used are French, German, Hindi, Korean, Arabic, Czech, Bengali, Vietnamese, Turkish, Telugu and Farsi.