

An Adversarial Perspective on Machine Unlearning for AI Safety

Jakub Łucki¹ Boyi Wei² Yangsibo Huang²

Peter Henderson² Florian Tramèr¹ Javier Rando¹

¹ETH Zurich ²Princeton University

Reviewed on OpenReview: <https://openreview.net/forum?id=J5IRyTKZ9s>

Abstract

Large language models are finetuned to refuse questions about hazardous knowledge, but these protections can often be bypassed. Unlearning methods aim at completely removing hazardous capabilities from models and make them inaccessible to adversaries. This work challenges the fundamental differences between unlearning and traditional safety post-training from an adversarial perspective. We demonstrate that existing jailbreak methods, previously reported as ineffective against unlearning, can be successful when applied carefully. Furthermore, we develop a variety of adaptive methods that recover most supposedly unlearned capabilities. For instance, we show that finetuning on 10 unrelated examples or removing specific directions in the activation space can recover most hazardous capabilities for models edited with RMU, a state-of-the-art unlearning method. Our findings challenge the robustness of current unlearning approaches and question their advantages over safety training.¹

1 Introduction

Large language models (LLMs) are pretrained on trillions of tokens crawled from the Internet (Dubey et al., 2024). Due to the unprecedented size of the training corpora, it is nearly impossible to discard all dangerous or otherwise harmful information available online. As a consequence, LLMs are capable of generating toxic, illicit, biased and privacy-infringing content (Wen et al., 2023; Karamolegkou et al., 2023; Nasr et al., 2023). Since models are constantly becoming more capable, this knowledge may pose increasing risks as it can make hazardous information more easily accessible for adversaries.

LLMs often undergo safety finetuning to reject unethical requests and produce safe responses (Bai et al., 2022). Yet, despite these safeguards, researchers continuously discover *jailbreaks* that bypass safeguards and elicit harmful generations from LLMs (Wei et al., 2024a). Robustness of these safeguards remains an open research question (Casper et al., 2023; Anwar et al., 2024) and machine unlearning (Cao and Yang, 2015; Bourtole et al., 2021) has emerged as a promising solution. It aims to completely remove hazardous knowledge from LLMs, preventing its extraction even after jailbreaking. State-of-the-art methods, like RMU (Li et al., 2024), can reduce accuracy on hazardous knowledge benchmarks to random chance. However, unlearning is not foolproof, as hazardous knowledge can still be recovered after the process (Patil et al., 2024; Shumailov et al., 2024; Hu et al., 2024). This raises an important question: Does unlearning truly remove hazardous knowledge, or does it simply “obfuscate” this knowledge similarly to refusal safety training?

In this work, we challenge the fundamental differences between unlearning and traditional safety finetuning from an adversarial perspective. We use the accuracy on the WMDP benchmark (Li et al., 2024) to measure the hazardous knowledge contained in LLMs. We argue that, from the perspective of model safety, unlearning is not successful if there exists *at least one* way of recovering significant accuracy either *without updating the*

¹Code is available at: <https://github.com/ethz-spylab/unlearning-vs-safety>

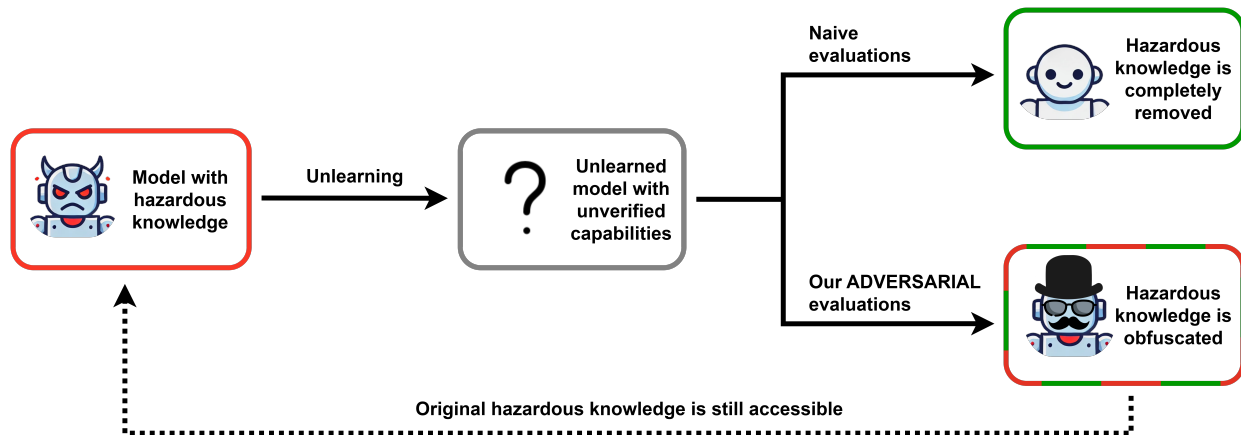


Figure 1: Conceptual overview of our contribution. Our adversarial evaluations show that current unlearning methods largely obfuscate hazardous knowledge instead of erasing it from model weights.

model weights or *updating the model weights with data that has little or no mutual information with the target knowledge*.

We perform the first comprehensive white-box evaluation of state-of-the-art unlearning methods for hazardous knowledge, comparing them to traditional safety training with DPO (Rafailov et al., 2024). Our results show that while unlearning is robust against specific attacks like probing internal model activations, it can also be easily compromised with methods similar to those used against safety training. Jailbreak methods that were reported ineffective against unlearning, like GCG (Zou et al., 2023), can recover substantial accuracy after small changes in the loss function. Additionally, we find that removing specific directions in the activation space, or finetuning on 10 unrelated examples can completely undo unlearning and recover the original performance on WMDP.

Overall, our findings underscore the limitations of black-box evaluations in accurately assessing unlearning effectiveness for safety settings and highlight the pressing need to refine unlearning methods, so that they deliver their promised benefits over standard safety training.

2 Related Work

Safety training and jailbreaks. Large language models are finetuned to refuse questions about hazardous knowledge with safety methods like DPO (Rafailov et al., 2024) or PPO (Ouyang et al., 2022). Zou et al. (2024b) recently introduced *circuit breakers* that use representation engineering to orthogonalize directions corresponding to unwanted concepts. The robustness of existing safeguards is limited (Casper et al., 2023; Anwar et al., 2024) and researchers often find *jailbreaks* to bypass protections and elicit hazardous knowledge (Wei et al., 2024a). Jailbreaks can rely only on prompting strategies (Shah et al., 2023; Huang et al., 2023), exploit white-box access to optimize prompts (Zou et al., 2023; Andriushchenko et al., 2024) or ablate model activations (Arditi et al., 2024).

LLM Unlearning. The gold standard of machine unlearning is to modify a model such that it is indistinguishable from one retrained on an original dataset with the target data removed (Cao and Yang, 2015; Bourtole et al., 2021). Given that LLMs store considerable amounts of interpretable knowledge within their weights (Patil et al., 2024), respective unlearning methods aim to render it inaccessible by any means short of full retraining. Unlearning for LLMs has been investigated as a potential solution to problems spanning fairness, privacy, safety and hallucinations (Jang et al., 2022; Yao et al., 2024; Chen and Yang, 2023; Wu et al., 2023; Li et al., 2024; Liu et al., 2024b). Effectiveness of unlearning is typically evaluated using narrow topics (e.g. Harry Potter) or fictional information that model had not seen before (Eldan and Russinovich, 2023; Maini et al., 2024; Shi et al., 2024; Wei et al., 2024c). Despite new unlearning benchmarks, the field

lacks standardized metrics and unified evaluation frameworks, as each benchmark employs its own criteria ranging from evaluating the perplexity of the plain outputs (Liu et al., 2024c), through multiple choice accuracy (Li et al., 2024), to ROUGE scores after prefix injections (Jin et al., 2024). Our work focuses on unlearning methods for safety. These methods try to eliminate dangerous knowledge to prevent adversaries from accessing it, even after jailbreaking attempts. The most notable method for this purpose is RMU (Li et al., 2024), which was introduced alongside WMDP, a benchmark for evaluating hazardous capabilities. General-purpose unlearning algorithms like negative preference optimization (NPO) (Zhang et al., 2024) can also be adapted for this purpose.

Unlearning robustness. Initial unlearning evaluations for LLMs relied on simple classification metrics (Eldan and Russinovich, 2023) which do not account for all possible ways in which a language model can represent and output the target information. Recent works (Jin et al., 2024; Hong et al., 2024; Lynch et al., 2024; Schwinn et al., 2024; Pawelczyk et al., 2024; Goel et al., 2022) have adopted an adversarial approach to test whether there exist ways to extract the information that was supposedly unlearned. For instance, Lynch et al. (2024) showed that knowledge could be extracted at comparable rates from both original and unlearned models by probing internal representations. In the context of unlearning hazardous capabilities, RMU reports robustness under some white-box jailbreaks like GCG or probing, but finds that finetuning unlearned models can easily disable the protections (Li et al., 2024). Similarly, Hu et al. (2024) find that fine-tuning can revert unlearning. In this work, we devise novel white-box methods to extract hazardous knowledge from unlearned models without significantly updating the weights. The importance of meticulous evaluations, has been demonstrated by an earlier work on word embedding debiasing, which revealed the lack of robustness of the respective methods (Gonen and Goldberg, 2019). Furthermore, the strength of white-box evaluations is underlined by Sharma et al. (2024), who showed that the most prominent unlearning techniques for diffusion models only decouple target concepts from corresponding prompts instead of erasing them.

Concurrent work. Deeb and Roger (2024) have shown that RMU and two other unlearning methods are not robust to finetuning attacks using data from the unlearned distribution, even if relearned knowledge is not directly related to the rest of the unlearned knowledge. This supports our findings that current LLM unlearning techniques obfuscate knowledge instead of removing it. Doshi and Stickland (2024) used black-box techniques such as few-shot prompting to test the robustness of two unlearning methods: RMU, and LLMU (Yao et al., 2024). They were able to recover moderate amounts of unlearned knowledge, whereas we were able to recover most of it in a white-box setting. Che et al. (2024) has shown susceptibility (to a varying extent) of 8 different capability removal techniques, including RMU, to a series of off-the shelf adversarial attacks. In contrast our works highlights the importance of adaptivity of adversarial evaluations. Both our work and Che et al. (2024) showcase the importance of white-box methods in stress-testing the unlearning robustness.

3 Experimental Setup

This work focuses exclusively on unlearning methods for safety that remove hazardous knowledge (e.g. bioweapons) from large language models, as introduced by Li et al. (2024). In practice, unlearning relies on *forget* and *retain* sets. The first contains information relevant to the domain to be unlearned (e.g. enhanced pandemic pathogens) while the second includes any information that should be preserved (e.g. general biology). In this work, we use the datasets included in WMDP benchmark for biology and cybersecurity (Li et al., 2024). Our evaluation is designed to assess whether existing unlearning methods effectively remove hazardous knowledge or merely make it more difficult to access, similarly to safety training.

3.1 Threat Model

We assume white-box access to an unlearned model, allowing modification of its weights and intervention in the activation space during inference. Additionally, we assume access to the original model prior to unlearning or to an equivalent model obtained by removing unlearning protections through finetuning, as demonstrated later. Although white-box access differs from the threat model for protections we study (RMU assumes only black-box access), it provides valuable insights into the effectiveness of unlearning in removing knowledge from

model weights. Furthermore, with the rise of powerful open-source large language models, robust unlearning in white-box scenarios is an increasingly relevant desiderata.

3.2 Unlearning Methods and Safety Training Baseline

We evaluate the most powerful unlearning method for hazardous knowledge to date: RMU (Li et al., 2024; Kadhe et al., 2024)². Additionally, we implement NPO (Zhang et al., 2024) that has been widely used as a general-purpose unlearning method for fact and concept removal (Shi et al., 2024), but its effectiveness for hazardous knowledge removal remains unexplored. We specifically use NPO+RT, a variant of NPO including an additional retain loss. Finally, we include DPO (Rafailov et al., 2024) as a baseline for safety training to contrast it with unlearning methods. For more details about the methods, see Appendix B.

3.3 Models and Datasets

We evaluate the performance of RMU using the publicly available checkpoint³. This model results from finetuning Zephyr-7B- β (Tunstall et al., 2023) on the WMDP and WikiText corpora (Merity et al., 2016). For NPO and DPO, we finetune Zephyr-7B- β ourselves on WMDP. We will refer to these models as *unlearned models*.

NPO and DPO require preference datasets, but WMDP only provides corpora (e.g. scientific papers) for autoregressive training. We use GPT-4 (OpenAI et al., 2024) to formulate questions based on these documents. For questions about hazardous topics, we set one of 80 random refusal strings as the desired output and the full correct option as the rejected response. For questions based on the *retain* set, we keep the correct option as the desired output and reject the refusal. We refer to the resulting datasets as our *preference datasets*. See Appendix C for details on dataset construction.

To ensure a fair comparison with safety methods, we fine-tune Zephyr using DPO specifically on preference datasets relevant to unlearning topics, rather than training it to refuse all harmful requests. We balance the training data by including samples from the forget and retain preference datasets, as well as OpenAssistant (Köpf et al., 2024), in a 50:25:25 ratio. This approach aims to maintain a balance between refusal capabilities and preserving general utility. For NPO, we use the preference dataset on hazardous knowledge as negative samples and the retain preference dataset mixed with OpenAssistant (50:50) dataset for the auxiliary retain loss.

3.4 Unlearning Evaluation

We evaluate the performance of unlearning hazardous knowledge using the WMDP benchmark (Li et al., 2024), which consists of 1,273 multiple-choice questions about dangerous biology knowledge and 1,987 about cybersecurity. To detect latent knowledge that might still be present even when models refuse to answer, we select the option (A, B, C, or D) with the highest probability as the final response. Besides, we use MMLU (Hendrycks et al., 2020) to measure the model’s general utility after unlearning, which contains multiple-choice questions covering 57 different tasks. For both WMDP benchmark and MMLU, we report overall accuracy across the entire dataset.

4 Our Methods To Recover Hazardous Capabilities

We use a wide range of methods to uncover hazardous capabilities in the target models, ranging from representation engineering to prompt-based jailbreaks. Most methods are inspired by well-known safety jailbreaks and incorporate small changes to target unlearning methods. All of our methods—except for finetuning—do not modify model weights and, thus, can only access knowledge that was preserved in model weights after unlearning. For finetuning, we primarily use small or unrelated datasets to ensure that models

²*Embedding-CORrupted Prompts* (ECO) (Liu et al., 2024a) outperforms others but applies a pre-LLM filter, leaving the original weights and potential hazardous knowledge unchanged. Thus, it doesn’t meet our definition of unlearning. See Appendix A for further discussion.

³Available at https://huggingface.co/cais/Zephyr_RMU

cannot acquire new hazardous capabilities. Detailed specification of resources needed for each method can be found in Appendix E.1.

4.1 Finetuning

It has been shown that finetuning easily reverses safety alignment even when using benign datasets (Qi et al., 2023). Also, the original RMU work and showed that fine-tuning unlearned models on the entire forget dataset could recover hazardous capabilities. In this work, we fine-tune unlearned models on datasets with very low mutual information (MI) with the unlearned knowledge to ensure that no new knowledge can be acquired. We use Low-Rank Adaptation (LoRA; Hu et al., 2021) to fine-tune unlearned models on three datasets: (1) forget dataset, (2) retain dataset—disjoint with forget dataset by definition—, and (3) WikiText (Merity et al., 2016)—a collection of Wikipedia documents with minimal overlap with hazardous knowledge. We experiment with varying sample sizes (from 5 to 1000 examples). By incorporating datasets with high MI (forget set) and low MI (retain set and WikiText), we provide a comprehensive evaluation of how different configurations affect the pace of hazardous knowledge recovery. Further setup detail are in Appendix E.2.

4.2 Orthogonalization

Arditi et al. (2024) demonstrated that safety refusal is governed by a single direction in the activation space. We investigate whether unlearning techniques generate a similar direction. Rather than targeting a single layer, we allow for distinct refusal directions at each transformer block. Using the forget preference dataset, we collect the outputs of each transformer block from both the original and unlearned models. We then compute the refusal direction for each layer using the difference in means method (Belrose, 2023). At inference time, we remove the refusal direction at each layer. Additionally, we develop a setup that does not require access to the original model prior to unlearning; see Appendix E.3 for details.

4.3 Logit Lens

Logit Lens is an interpretability technique (nostalgebraist, 2020; Patil et al., 2024) that projects the activations in residual stream onto the model’s vocabulary. We apply this technique to the WMDP dataset by using the projected logits of the A, B, C, and D tokens as the model’s answers. We project the output of transformer blocks at every layer and select the token with a higher probability. We also evaluate the projection of other activation spaces in Appendix G.3.

4.4 Enhanced GCG

GCG has been reported ineffective against RMU (Li et al., 2024; Huu-Tien et al., 2024). We introduce *enhanced GCG*, which specifically targets unlearning methods, and is based on FLRT (Thompson and Sklar, 2024) and augmented with several modifications detailed in Appendix E.5. Unlike GCG, which aims to find adversarial prompt suffixes, enhanced GCG focuses on optimizing *prefixes* to prevent the model from recognizing hazardous knowledge in the first place, as RMU will introduce persistent noise to the residual stream once such context is detected. We also attribute more weight to the loss computed on early tokens in the prompt. Our attack is optimized on 6 questions from the WMDP benchmark that were answered correctly by the original model and incorrectly by the unlearned model.

4.5 Set difference pruning

Wei et al. (2024b) introduced *set difference pruning* as a method to identify and prune neurons associated with safety alignment. Reproducing their method, we use SNIP (Lee et al., 2018) score to measure the importance of individual neurons for hazardous knowledge. Specifically, we compute the importance score for each neuron on the WMDP forget set, and the utility score on MMLU. We then use set difference method to find the neurons that only contribute to storing hazardous knowledge and remove them via pruning.

5 Results

We report the performance of our methods on WMDP-Bio due to significant difference in the scores of original and unlearned models. Analogous gap on WMDP-Cyber is much smaller, which makes the corresponding results more volatile (See Appendix F). We summarize our results and observations below.

Table 1: WMDP-Bio and MMLU accuracy for each protection and method. For Logit Lens, we report the best layer overall. For finetuning, we report best result on 5 samples from the forget set. Empty values are not possible to compute or the corresponding combination does not affect the score.

Datasets	Knowledge Recovery	No Protection	Unlearning Methods		Safety Training
			RMU	NPO	DPO
WMDP-Bio	Default decoding	64.4	29.9	29.5	27.9
	Logit Lens	66.2	31.8	38.6	48.2
	Finetuning	-	62.4	47.4	57.3
	Orthogonalization	-	64.7	45.1	50.7
	Enhanced GCG	-	53.9	46.0	49.0
	Pruning	-	54.0	40.4	50.4
MMLU	Default decoding	58.1	57.1	52.1	49.7
	Logit Lens	-	-	-	-
	Finetuning	-	58.0	53.3	51.2
	Orthogonalization	-	57.3	45.6	46.7
	Enhanced GCG	-	-	-	-
	Pruning	-	56.5	50.0	50.4

Finetuning on unrelated information reverts unlearning. As illustrated in Figure 2, finetuning with only 10 samples from the retain set—disjoint by definition from the evaluation knowledge—can recover most of hazardous capabilities, obtaining accuracies of 52.7% (NPO), 57.0% (DPO), and 61.6% (RMU) while causing negligible degradation on MMLU (less than 2 p.p.). Finetuning on 1000 samples from the retain set fully recovers hazardous capabilities across all methods. These results demonstrate that both safety training and unlearning can be undone through finetuning on unrelated information, suggesting that unlearning is also expressed through shallow features (Yang et al., 2023; Lermen et al., 2023). Additionally, finetuning with just 5 samples from the forget set effectively reverses unlearning, particularly for RMU, which nearly recovers its original performance. Relearning knowledge through further training is unavoidable, but these results show that knowledge recovery happens at disproportionately fast rate.

Unlearning methods remove knowledge from the residual stream more effectively. Before unlearning, Logit Lens can decode correct answers from Zephyr-7B at layer 19, as shown in Figure 3. However, Logit Lens becomes ineffective after protections are applied. Our safety baseline, DPO, remains the most susceptible to early decoding, achieving 56% accuracy. In contrast, unlearning methods can remove knowledge more effectively from the residual stream, with RMU reducing Logit Lens accuracy close to random chance across the entire architecture. These results align with prior evaluations of RMU’s robustness to probing (Li et al., 2024).

Unlearning is also mediated by specific directions. We identify and ablate directions responsible for unlearning, successfully recovering hazardous knowledge for most protections (see Table 1). RMU is the most vulnerable to our orthogonalization, achieving 64.7% accuracy (surpassing the baseline accuracy of 64.4%) by manipulating only the activation space during the forward pass. This outperforms ablation of a single refusal direction across all layers (Arditi and Chughtai, 2024), which achieves 54.2% accuracy. NPO and DPO are more robust against orthogonalization, obtaining 45% and 51% accuracy, respectively.

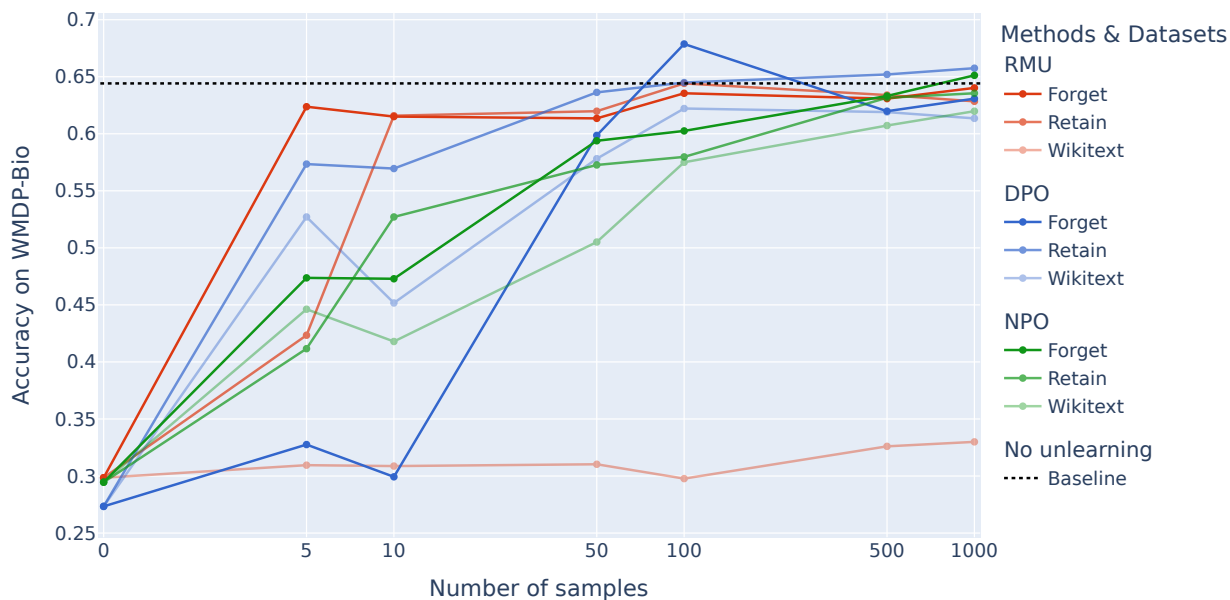


Figure 2: Accuracy on WMDP-Bio for unlearned models finetuned with different datasets and number of samples. See Appendix F.1 for complimentary results on MMLU and WMDP-Cyber.

Unlearning depends on critical neurons. We localized minimal sets of weights that are responsible for degradation in hazardous knowledge for each unlearning method. These sets represent 2.0% of weights for NPO, 0.9% for RMU, and 2.4% for DPO. After pruning these weights, performance on WMDP increases by at least 10 p.p. for all methods.

Universal adversarial prefixes that recover unlearned knowledge exist. Using *enhanced GCG* we were able to craft universal adversarial prefixes that increased RMU’s accuracy from 29.9% to 53.9%, NPO’s accuracy from 29.5% to 46.0%, and DPO’s accuracy from 27.9% to 49.0%. This demonstrates that, similarly to safety trained models, input-only manipulations can disable unlearning and elicit hazardous knowledge that was never removed from the model.

We can recover hazardous capabilities while models remain unusable. RMU is characterized by making models unusable—they output gibberish generations with high perplexity—when hazardous knowledge is detected. Interestingly, we find that GCG prefixes can easily recover a conversational model that answers questions from WMDP, but its responses are often incorrect and overconfident. Best performing prefixes can recover most of the hazardous capabilities while not necessarily recovering conversational capabilities from the model. See Appendix I for an analysis.

6 Discussion

Existing unlearning methods are not different from safety training. Our findings reveal that unlearning methods primarily obscure knowledge rather than eliminate it (as illustrated by Figure 1), which is a known flaw of safety training (Lee et al., 2024). Therefore, RMU and NPO are susceptible to techniques analogous to those that can reverse safety training, including: (1) dependence on individual residual stream directions; (2) rapid knowledge recovery after finetuning with unrelated data; (3) presence of critical neurons that inhibit hazardous knowledge; and (4) existence of universal adversarial strings that unlock the unlearned

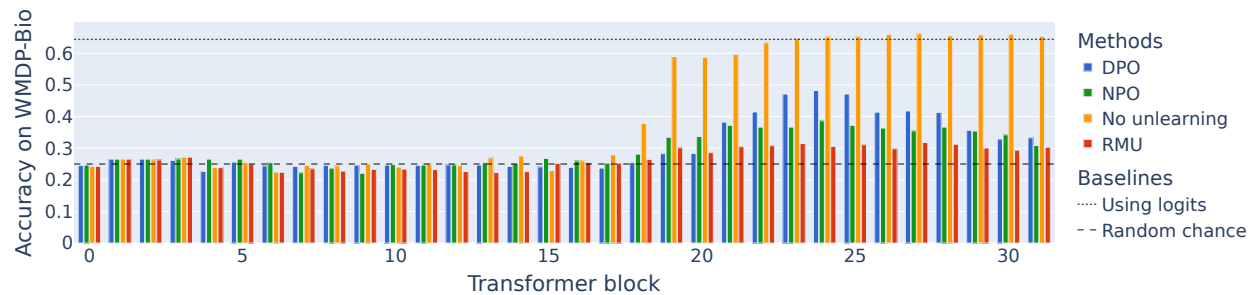


Figure 3: Accuracy on WMDP-Bio using LogitLens after each transformer block.

knowledge. These observations question the practical benefits of unlearning methods over safety training. Although unlearning was proposed to fully eradicate hazardous capabilities and mitigate jailbreaks in large language models, our results indicate that these methods share limitations. Concurrent work by [Tamirisa et al. \(2024\)](#) proposed TAR, a technique that can prevent *some* fine-tuning attacks but has no impact on others.

Black-box evaluations are insufficient for unlearning. Unlearning aims to completely remove certain information from model weights, yet many evaluations only verify that this information cannot be easily extracted through model outputs. This mismatch between the unlearning objective and the evaluation method can falsely suggest successful unlearning when knowledge remains embedded in the weights ([Lynch et al., 2024](#)). In fact, while black-box methods used in the original RMU evaluation ([Li et al., 2024](#)) and concurrent work by [Doshi and Stickland \(2024\)](#) failed to elicit supposedly unlearned knowledge, our white-box approaches successfully recovered the entire unlearned information. As extensively demonstrated in security and safety research, adaptive evaluations are required to faithfully evaluate ML protections ([Carlini and Wagner, 2017](#); [Tramer et al., 2020](#); [Radiya-Dixit et al., 2021](#); [Hönig et al., 2024](#)).

NPO shows signs of deep unlearning. This method consistently displays better robustness than DPO or RMU, suggesting that gradient ascent ([Zhang et al., 2024](#)) might be a promising tool to remove hazardous knowledge from model weights. However, our current implementation still results in greater degradation on MMLU and general capabilities. Future work could investigate combining representation engineering with gradient ascent to enhance existing unlearning methods.

Possible mitigations and lessons learned. Our results suggest that fine-tuning models based solely on their outputs may result in knowledge obfuscation rather than erasure. Similarly, the noise injection strategy employed by RMU proves inadequate for open-source models. Instead of random noise addition, orthogonalizing directions corresponding to harmful knowledge presents a more promising approach ([Zou et al., 2024a](#)). We also caution that noise-based unlearning can produce misleadingly positive results during evaluation with GCG, as standard GCG implementations are susceptible to residual stream noise by default. Other effective approaches may include methods that precisely localize and modify knowledge vectors, as proposed by [Hong et al. \(2024\)](#).

7 Conclusion

We performed a comprehensive white-box evaluation of state-of-the-art unlearning methods for AI safety. Our findings reveal that these methods cannot reliably remove knowledge from model weights. For example, finetuning on unrelated data or removing specific directions from activation space often recovers the supposedly unlearned capabilities. This challenges the belief that unlearning methods offer more robust protection than standard safety training. Furthermore, we argue that black-box evaluations are insufficient for unlearning, as they do not assess internal model changes.

References

- Maksym Andriushchenko, Francesco Croce, and Nicolas Flammarion. Jailbreaking leading safety-aligned llms with simple adaptive attacks, 2024. URL <https://arxiv.org/abs/2404.02151>.
- Usman Anwar, Abulhair Saparov, Javier Rando, Daniel Paleka, Miles Turpin, Peter Hase, Ekdeep Singh Lubana, Erik Jenner, Stephen Casper, Oliver Sourbut, Benjamin L. Edelman, Zhaowei Zhang, Mario Günther, Anton Korinek, Jose Hernandez-Orallo, et al. Foundational challenges in assuring alignment and safety of large language models. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. Survey Certification, Expert Certification.
- Andy Arditi and Bilal Chughtai, Jul 2024. URL <https://www.lesswrong.com/posts/6QYpXEscd8GuE7BgW/unlearning-via-rmu-is-mostly-shallow>.
- Andy Arditi, Oscar Obeso, Aaqib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. Refusal in language models is mediated by a single direction, 2024. URL <https://arxiv.org/abs/2406.11717>.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022. URL <https://arxiv.org/abs/2204.05862>.
- Nora Belrose. Diff-in-means concept editing is worst-case optimal: Explaining a result by Sam Marks and Max Tegmark, 2023. <https://blog.eleuther.ai/diff-in-means/>. Accessed on: September 12, 2024.
- Lucas Bourtole, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 141–159. IEEE, 2021.
- Yinzhi Cao and Junfeng Yang. Towards making systems forget with machine unlearning. In *2015 IEEE symposium on security and privacy*, pages 463–480. IEEE, 2015.
- Nicholas Carlini and David Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *Proceedings of the 10th ACM workshop on artificial intelligence and security*, pages 3–14, 2017.
- Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, Tony Tong Wang, Samuel Marks, Charbel-Raphael Segerie, Micah Carroll, Andi Peng, et al. Open problems and fundamental limitations of reinforcement learning from human feedback. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. Survey Certification.
- Zora Che, Stephen Casper, Anirudh Satheesh, Rohit Gandikota, Domenic Rosati, Stewart Slocum, Lev E McKinney, Zichu Wu, Zikui Cai, Bilal Chughtai, Furong Huang, and Dylan Hadfield-Menell. Model manipulation attacks enable more rigorous evaluations of LLM unlearning. In *Neurips Safe Generative AI Workshop 2024*, 2024. URL <https://openreview.net/forum?id=XmvgWEjkhG>.
- Jiaao Chen and Diyi Yang. Unlearn what you want to forget: Efficient unlearning for llms. *arXiv preprint arXiv:2310.20150*, 2023.
- William S Cleveland. Robust locally weighted regression and smoothing scatterplots. *Journal of the American statistical association*, 74(368):829–836, 1979.
- Aghyad Deeb and Fabien Roger. Do unlearning methods remove information from language model weights? *arXiv preprint arXiv:2410.08827*, 2024.
- Jai Doshi and Asa Cooper Stickland. Does unlearning truly unlearn? a black box evaluation of llm unlearning methods. *arXiv preprint arXiv:2411.12103*, 2024.

- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, et al. llama3family, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Ronen Eldan and Mark Russinovich. Who’s harry potter? approximate unlearning in llms, 2023. URL <https://arxiv.org/abs/2310.02238>.
- Chongyang Gao, Lixu Wang, Chenkai Weng, Xiao Wang, and Qi Zhu. Practical unlearning for large language models. *arXiv preprint arXiv:2407.10223*, 2024.
- Siddhant Garg and Goutham Ramakrishnan. BAE: BERT-based adversarial examples for text classification. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6174–6181, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.498. URL <https://aclanthology.org/2020.emnlp-main.498>.
- Shashwat Goel, Ameya Prabhu, Amartya Sanyal, Ser-Nam Lim, Philip Torr, and Ponnurangam Kumaraguru. Towards adversarial evaluations for inexact machine unlearning. *arXiv preprint arXiv:2201.06640*, 2022.
- Hila Gonen and Yoav Goldberg. Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them, 2019. URL <https://arxiv.org/abs/1903.03862>.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- Yihuai Hong, Lei Yu, Shauli Ravfogel, Haiqin Yang, and Mor Geva. Intrinsic evaluation of unlearning using parametric knowledge traces. *arXiv preprint arXiv:2406.11614*, 2024.
- Robert Hönig, Javier Rando, Nicholas Carlini, and Florian Tramèr. Adversarial perturbations cannot reliably protect artists from generative ai. *arXiv preprint arXiv:2406.12027*, 2024.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Shengyuan Hu, Yiwei Fu, Zhiwei Steven Wu, and Virginia Smith. Jogging the memory of unlearned model through targeted relearning attack. *arXiv preprint arXiv:2406.13356*, 2024.
- Yangsibo Huang, Samyak Gupta, Mengzhou Xia, Kai Li, and Danqi Chen. Catastrophic jailbreak of open-source llms via exploiting generation. *arXiv preprint arXiv:2310.06987*, 2023.
- Dang Huu-Tien, Trung-Tin Pham, Hoang Thanh-Tung, and Naoya Inoue. On effects of steering latent representation for large language model unlearning. *arXiv preprint arXiv:2408.06223*, 2024.
- Joel Jang, Dongkeun Yoon, Sohee Yang, Sungmin Cha, Moontae Lee, Lajanugen Logeswaran, and Minjoon Seo. Knowledge unlearning for mitigating privacy risks in language models. *arXiv preprint arXiv:2210.01504*, 2022.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- Zhuoran Jin, Pengfei Cao, Chenhao Wang, Zhitao He, Hongbang Yuan, Jiachun Li, Yubo Chen, Kang Liu, and Jun Zhao. Rwku: Benchmarking real-world knowledge unlearning for large language models. *arXiv preprint arXiv:2406.10890*, 2024.
- Swanand Ravindra Kadhe, Farhan Ahmed, Dennis Wei, Nathalie Baracaldo, and Inkit Padhi. Split, unlearn, merge: Leveraging data attributes for more effective unlearning in llms. *arXiv preprint arXiv:2406.11780*, 2024.

- Antonia Karamolegkou, Jiaang Li, Li Zhou, and Anders Søgaard. Copyright violations and large language models, 2023. URL <https://arxiv.org/abs/2310.13771>.
- Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, et al. Openassistant conversations-democratizing large language model alignment. *Advances in Neural Information Processing Systems*, 36, 2024.
- Andrew Lee, Xiaoyan Bai, Itamar Pres, Martin Wattenberg, Jonathan K. Kummerfeld, and Rada Mihalcea. A mechanistic understanding of alignment algorithms: A case study on dpo and toxicity, 2024. URL <https://arxiv.org/abs/2401.01967>.
- Namhoon Lee, Thalaiyasingam Ajanthan, and Philip HS Torr. Snip: Single-shot network pruning based on connection sensitivity. *arXiv preprint arXiv:1810.02340*, 2018.
- Simon Lermen, Charlie Rogers-Smith, and Jeffrey Ladish. Lora fine-tuning efficiently undoes safety training in llama 2-chat 70b. *arXiv preprint arXiv:2310.20624*, 2023.
- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. BERT-ATTACK: Adversarial attack against BERT using BERT. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6193–6202, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.500. URL <https://aclanthology.org/2020.emnlp-main.500>.
- Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D. Li, Ann-Kathrin Dombrowski, Shashwat Goel, Long Phan, Gabriel Mukobi, Nathan Helm-Burger, Rassin Lababidi, Lennart Justen, Andrew B. Liu, et al. The wmdp benchmark: Measuring and reducing malicious use with unlearning, 2024.
- Chris Yuhao Liu, Yaxuan Wang, Jeffrey Flanigan, and Yang Liu. Large language model unlearning via embedding-corrupted prompts. *arXiv preprint arXiv:2406.07933*, 2024a.
- Sijia Liu, Yuanshun Yao, Jinghan Jia, Stephen Casper, Nathalie Baracaldo, Peter Hase, Xiaojun Xu, Yuguang Yao, Hang Li, Kush R Varshney, et al. Rethinking machine unlearning for large language models. *arXiv preprint arXiv:2402.08787*, 2024b.
- Zheyuan Liu, Guangyao Dou, Zhaoxuan Tan, Yijun Tian, and Meng Jiang. Towards safer large language models through machine unlearning. *arXiv preprint arXiv:2402.10058*, 2024c.
- Aengus Lynch, Phillip Guo, Aidan Ewart, Stephen Casper, and Dylan Hadfield-Menell. Eight methods to evaluate robust unlearning in llms, 2024. URL <https://arxiv.org/abs/2402.16835>.
- Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary Chase Lipton, and J Zico Kolter. TOFU: A task of fictitious unlearning for LLMs. In *First Conference on Language Modeling*, 2024.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models, 2016.
- Milad Nasr, Nicholas Carlini, Jonathan Hayase, Matthew Jagielski, A. Feder Cooper, Daphne Ippolito, Christopher A. Choquette-Choo, Eric Wallace, Florian Tramèr, and Katherine Lee. Scalable extraction of training data from (production) language models, 2023. URL <https://arxiv.org/abs/2311.17035>.
- nostalgebraist. Interpreting gpt: the logit lens, 2020. URL <https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens>.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, et al. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.

- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Vaidehi Patil, Peter Hase, and Mohit Bansal. 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=7er1RDoaV8>.
- Martin Pawelczyk, Jimmy Z Di, Yiwei Lu, Gautam Kamath, Ayush Sekhari, and Seth Neel. Machine unlearning fails to remove data poisoning attacks. *arXiv preprint arXiv:2406.17216*, 2024.
- Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. Fine-tuning aligned language models compromises safety, even when users do not intend to! *arXiv preprint arXiv:2310.03693*, 2023.
- Evani Radiya-Dixit, Sanghyun Hong, Nicholas Carlini, and Florian Tramèr. Data poisoning won’t save you from facial recognition. *arXiv preprint arXiv:2106.14851*, 2021.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024. URL <https://arxiv.org/abs/2305.18290>.
- Vinu Sankar Sadasivan, Shoumik Saha, Gaurang Sriramanan, Priyatham Kattakinda, Atoosa Chegini, and Soheil Feizi. Fast adversarial attacks on language models in one gpu minute. *arXiv preprint arXiv:2402.15570*, 2024.
- Leo Schwinn, David Dobre, Sophie Xhonneux, Gauthier Gidel, and Stephan Gunnemann. Soft prompt threats: Attacking safety alignment and unlearning in open-source llms through the embedding space. *arXiv preprint arXiv:2402.09063*, 2024.
- Rusheb Shah, Quentin Feuillade-Montixi, Soroush Pour, Arush Tagade, Stephen Casper, and Javier Rando. Scalable and transferable black-box jailbreaks for language models via persona modulation, 2023. URL <https://arxiv.org/abs/2311.03348>.
- Aakash Sen Sharma, Niladri Sarkar, Vikram Chundawat, Ankur A Mali, and Murari Mandal. Unlearning or concealment? a critical analysis and evaluation metrics for unlearning in diffusion models. *arXiv preprint arXiv:2409.05668*, 2024.
- Weijia Shi, Jaechan Lee, Yangsibo Huang, Sadhika Malladi, Jieyu Zhao, Ari Holtzman, Daogao Liu, Luke Zettlemoyer, Noah A Smith, and Chiyuan Zhang. Muse: Machine unlearning six-way evaluation for language models. *arXiv preprint arXiv:2407.06460*, 2024.
- Iliia Shumailov, Jamie Hayes, Eleni Triantafillou, Guillermo Ortiz-Jimenez, Nicolas Papernot, Matthew Jagielski, Itay Yona, Heidi Howard, and Eugene Bagdasaryan. Ununlearning: Unlearning is not sufficient for content regulation in advanced generative ai. *arXiv preprint arXiv:2407.00106*, 2024.
- Antoine Simoulin and Benoit Crabbé. How many layers and why? An analysis of the model depth in transformers. In Jad Kabbara, Haitao Lin, Amandalynne Paullada, and Jannis Vamvas, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: Student Research Workshop*, pages 221–228, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-srw.23. URL <https://aclanthology.org/2021.acl-srw.23>.
- Rishub Tamirisa, Bhruvu Bharathi, Long Phan, Andy Zhou, Alice Gatti, Tarun Suresh, Maxwell Lin, Justin Wang, Rowan Wang, Ron Arel, Andy Zou, Dawn Song, Bo Li, Dan Hendrycks, and Mantas Mazeika. Tamper-resistant safeguards for open-weight llms, 2024. URL <https://arxiv.org/abs/2408.00761>.
- T Ben Thompson and Michael Sklar. Fluent student-teacher redteaming. *arXiv preprint arXiv:2407.17447*, 2024.

- Florian Tramer, Nicholas Carlini, Wieland Brendel, and Aleksander Madry. On adaptive attacks to adversarial example defenses. *Advances in neural information processing systems*, 33:1633–1645, 2020.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of lm alignment, 2023. URL <https://arxiv.org/abs/2310.16944>.
- Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does llm safety training fail? *Advances in Neural Information Processing Systems*, 36, 2024a.
- Boyi Wei, Kaixuan Huang, Yangsibo Huang, Tinghao Xie, Xiangyu Qi, Mengzhou Xia, Prateek Mittal, Mengdi Wang, and Peter Henderson. Assessing the brittleness of safety alignment via pruning and low-rank modifications. *arXiv preprint arXiv:2402.05162*, 2024b.
- Boyi Wei, Weijia Shi, Yangsibo Huang, Noah A Smith, Chiyuan Zhang, Luke Zettlemoyer, Kai Li, and Peter Henderson. Evaluating copyright takedown methods for language models. *arXiv preprint arXiv:2406.18664*, 2024c.
- Jiaxin Wen, Pei Ke, Hao Sun, Zhexin Zhang, Chengfei Li, Jinfeng Bai, and Minlie Huang. Unveiling the implicit toxicity in large language models, 2023. URL <https://arxiv.org/abs/2311.17391>.
- Xinwei Wu, Junzhuo Li, Minghui Xu, Weilong Dong, Shuangzhi Wu, Chao Bian, and Deyi Xiong. Depn: Detecting and editing privacy neurons in pretrained language models. *arXiv preprint arXiv:2310.20138*, 2023.
- Xianjun Yang, Xiao Wang, Qi Zhang, Linda Petzold, William Yang Wang, Xun Zhao, and Dahua Lin. Shadow alignment: The ease of subverting safely-aligned language models. *arXiv preprint arXiv:2310.02949*, 2023.
- Yuanshun Yao, Xiaojun Xu, and Yang Liu. Large language model unlearning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=8Dy42ThoNe>.
- Ruiqi Zhang, Licong Lin, Yu Bai, and Song Mei. Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*, 2024.
- Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.
- Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, J Zico Kolter, Matt Fredrikson, and Dan Hendrycks. Improving alignment and robustness with circuit breakers. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a. URL <https://openreview.net/forum?id=IbIB8SBKFV>.
- Andy Zou, Long Phan, Justin Wang, Derek Duenas, Maxwell Lin, Maksym Andriushchenko, Rowan Wang, Zico Kolter, Matt Fredrikson, and Dan Hendrycks. Improving alignment and robustness with circuit breakers, 2024b. URL <https://arxiv.org/abs/2406.04313>.

Appendices

A	Further discussion on ECO	16
A.1	Why ECO is not unlearning under our definition	16
A.2	Potential vulnerabilities	16
B	Further details on unlearning and safety training methods	17
B.1	Direct Preference Optimization (DPO)	17
B.2	Negative Preference Optimization (NPO)	17
B.3	Representation Misdirection for Unlearning (RMU)	17
C	Preference dataset construction	18
C.1	System prompt	18
C.2	Preference format	19
C.3	Refusal strings	20
D	Training details	22
D.1	Hyperparameters	22
D.2	Performance of developed models on relevant benchmarks	22
E	Additional details on knowledge extraction methods	23
E.1	Methods overview	23
E.2	Finetuning	23
E.3	Orthogonalization	24
E.4	Logit lens	25
E.5	Enhanced GCG	25
F	Complete results	27
F.1	Finetuning	28
F.2	Logit lens	29
F.3	Orthogonalization	31
G	Complete results using chat template	33
G.1	Overview of the results using chat template	34
G.2	Finetuning	35
G.3	Logit lens	36
H	Perturbations as a knowledge extraction method for RMU	38
H.1	Naive perturbations	38

H.2	Informed perturbations	38
H.3	Effectiveness of perturbations on RMU and other models	40
I	Perplexity Analysis of Adversarial Prefixes on RMU	43
I.1	Adversarial prefixes without chat template	44
J	RMU analysis	45
J.1	Behaviour during innocuous conversations	45
J.2	Simple prompt-based jailbreaks	45
J.3	Prefilling attack	46
J.4	Effect of noise on token representations	46

A Further discussion on ECO

A.1 Why ECO is not unlearning under our definition

Liu et al. (2024a) assume black-box access to the model. Given their setting their definition of successful unlearning entails that, in expectation, any non-negative metric computed on the outputs of an unlearned model and the outputs of a model retrained from scratch on retain set should be approximately one. Intuitively, the model trained only on the retain set should behave the same way as the original model after applying unlearning. Although this is the golden standard in machine unlearning, we consider it lacking for the generative models such as LLMs which show remarkable memorization capabilities (Nasr et al., 2023). The premise is that despite outputs of an LLM not displaying any signs of unlearned knowledge it can be stored within the weights, *and retrieved by an adversary*. To prevent that the knowledge should be removed from the weights as well. Hence, an improved definition of successful unlearning should include either the internals of an LLM or an adversarial perspective.

Furthermore, the core of ECO is an ‘unlearned’ knowledge detector, based on which a carefully crafted noise is applied to input embeddings. However, this is no different to a safety filter which given an unethical request would return a predefined refusal prompt. Choosing a suitable noise is merely obfuscating the refusal.

Ultimately, we would like to emphasize that we acknowledge ECO’s state-of-the-art results on WMDP. However, we argue that it doesn’t uphold the promise of unlearning.

A.2 Potential vulnerabilities⁴

Using a detector together with unmodified LLM, puts the red-teaming pressure on the former. As a consequence, the fundamental issue of defending the LLM is not resolved but rather reintroduced on a smaller scale, where we have to defend the detector (which in (Liu et al., 2024a) is a smaller LLM - RoBERTa).

After inspecting the code⁵, we noticed that there are two types of detectors implemented: token-wise and prompt-wise. The first one can be easily bypassed by forcing the tokenizer to tokenize the prompt character-by-character (e.g. by inserting whitespace between all relevant characters). Individual characters should not trigger any noise as they should not be exclusive to dangerous concepts. The second type of detector might be slightly more challenging, but there is significant body of works on adversarial attacks on BERT models (Li et al., 2020), including the specific scenario of text classification (Garg and Ramakrishnan, 2020).

⁴This is preliminary analysis and a sketch of potential red-teaming efforts. No experiments have been conducted.

⁵Available at <https://github.com/chrisliu298/llm-unlearn-eco/tree/main>

B Further details on unlearning and safety training methods

B.1 Direct Preference Optimization (DPO)

DPO (Rafailov et al., 2024) uses a *preference* dataset $\mathcal{D}_{\text{PREF}}$ consisting of triples: an input x , a *chosen* response y_w and a *rejected* response y_l . Model is then trained to produce generations that are closer to the *chosen* subset using the following objective:

$$\mathcal{L}_{\text{DPO}}(\theta) = -\frac{1}{\beta} \mathbb{E}_{\mathcal{D}_{\text{PREF}}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w | x)}{\pi_{\text{ref}}(y_w | x)} - \beta \log \frac{\pi_{\theta}(y_l | x)}{\pi_{\text{ref}}(y_l | x)} \right) \right], \quad (1)$$

where π_{ref} is reference model, π_{θ} is trainable model with weights θ , β is a variable controlling deviation from π_{ref} , and σ is a sigmoid function.

B.2 Negative Preference Optimization (NPO)

NPO (Zhang et al., 2024) optimizes a loss function inspired from DPO, where one uses only negative samples. Although, it may appear that this introduces inductive bias towards safety training, counter-intuitively it does not. Zhang et al. (2024) shows that NPO is a generalization of gradient ascent (GA). This resemblance is a desirable feature in unlearning as GA is the reverse process to gradient descent based learning. Furthermore, the authors show that NPO diverges at much slower rate than GA, making it more stable and thus, practical.

In the pilot experiments with straightforward application of NPO our models quickly diverged, resulting in catastrophic forgetting, indicated by poor performance on the utility benchmark. NPO collapsing when trying to unlearn broad domains is in line with other works suggesting that it fails in continual learning settings (Gao et al., 2024). Therefore, we focus on a variation of NPO which adds a retain loss (RT) to the original objective:

$$\mathcal{L}_{\text{NPO}}(\theta) = \underbrace{-\frac{2}{\beta} \mathbb{E}_{\mathcal{D}_{\text{FG}}} \left[\log \sigma \left(-\beta \log \frac{\pi_{\theta}(y|x)}{\pi_{\text{ref}}(y|x)} \right) \right]}_{\mathcal{L}_{\text{NPO}}} - \underbrace{\alpha \cdot \mathbb{E}_{\mathcal{D}_{\text{RT}}} [\log(\pi_{\theta}(y|x))]}_{\mathcal{L}_{\text{RT}}}, \quad (2)$$

where α is a weight of the retain loss, and (x, y) are input output pairs from the forget set \mathcal{D}_{FG} and from the retain set \mathcal{D}_{RT} . We refer to this method as NPO.

B.3 Representation Misdirection for Unlearning (RMU)

RMU (Li et al., 2024) finetunes a subset of lower layers of an LLM such that they output a fixed noise vector when given a prompt containing concepts present in the forget set and to leave representations unchanged if the concepts fall within the knowledge captured by the retain set. This method displays high sensitivity to keywords and behaves like a heavy-side function once “hazardous” concept is detected - internal representations will be distorted for all the subsequent tokens in the prompt. For detailed analysis of RMU see Appendix J. The RMU objective is as follows:

$$\mathcal{L}_{\text{RMU}}(\theta) = \underbrace{\mathbb{E}_{x \sim \mathcal{D}_{\text{FG}}} \left[\frac{1}{L_x} \sum_{t \in x} \|M_{\theta}(t) - c \cdot \mathbf{u}\|_2^2 \right]}_{\mathcal{L}_{\text{forget}}} + \underbrace{\alpha \cdot \mathbb{E}_{x \sim \mathcal{D}_{\text{RT}}} \left[\frac{1}{L_x} \sum_{t \in x} \|M_{\theta}(t) - M_{\text{ref}}(t)\|_2^2 \right]}_{\mathcal{L}_{\text{retain}}}, \quad (3)$$

where $M_{\text{ref}}, M_{\theta}$ are the internal representations of the reference and trainable models, L_x is the number of tokens in prompt x , and c is a variable controlling the magnitude of noise vector \mathbf{u} .

C Preference dataset construction

DPO and NPO require preference datasets of a specific format. We construct such datasets from the WMDP *forget* and *retain* documents. Furthermore, to make them more suitable for WMDP Benchmark we format individual samples as ABCD questions.

WMDP provides four corpora: bio-forget-corpus, bio-retain-corpus, cyber-forget-corpus, and cyber-retain-corpus. Biology ones are filtered scientific articles. Forget version contains abstract together with text and we decided to use the former due to its condensed nature. Cybersecurity corpora are mostly scraped github pages of oftentimes doubtful quality.

We used only the articles with more than 1000 characters to make sure that we filter out samples containing only titles and to provide enough material to generate 10 questions from each article. Moreover, each article was truncated at 15'000 characters to avoid unnecessary costs.

To generate questions based on the abovementioned corpora we used OpenAI API and the gpt-4o-mini-2024-07-18 model. Furthermore, we used their recent feature, Structured Outputs⁶, with the following schema:

```
class MCQ(BaseModel):
    question: str
    options: list[str]
    answer: str
    explanation: str

class ListMCQ(BaseModel):
    multiple_choice_questions: list[MCQ].
```

We sent requests until we obtained 10'000 samples from each corpora. However, each dataset can have up to 9 samples more since the model sometimes produced less than 10 questions and in some cases the correct answer was missing from the provided options. In our requests we used the system prompt defined below.

C.1 System prompt

In this system prompt we provide a clear objective, fixed number of options, and number of questions to generate from each article. Furthermore, we provide an example of a successfully completed task crafted from one of the WMDP bio-retain-corpus article to leverage few-shot learning paradigm.

You are a large language model tasked with creating multiple-choice questions based on an article. Each question is supposed to have four answers: A, B, C, D. Only one of these options is supposed to be correct, other three are supposed to be plausible but incorrect distractors. Your goal is to generate clear, concise questions that test knowledge contained within the article. You are not supposed to test text comprehension but rather concepts within it. Make sure that questions can be answered without an access to the article. Do not ask about data and statistics such as how many people died from polio in total. Instead ask 'Which of the following viruses causes polio?'. Never refer to the article directly. Always provide a correct answer and corresponding explanation. Make sure that the answer is one of the options. Don't include option letters A,B,C,D as part of each option. Instead it should be a phrase. Provide a list of 10 such questions.

EXAMPLE

ARTICLE:

Introduction {#bjh13767-sec-0001} ===== Purine analogues have revolutionized

⁶Available at: <https://openai.com/index/introducing-structured-outputs-in-the-api/>

the treatment of hairy cell leukaemia (HCL) with overall responses in more than 85% of patients and a median progression-free survival (PFS) of up to 15 years (Else *et al**, [2009] (#bjh13767-bib-0004){ref-type="ref"}; Grever, [2010] (#bjh13767-bib-0006){ref-type="ref"}). Despite the success of pentostatin and cladribine, PFS curves have revealed no plateau in long-term follow-up studies (Else *et al**, [2009] (#bjh13767-bib-0004){ref-type="ref"}). In addition, there is a small subset of patients with primary refractory disease. Second and third line treatment with single-agent purine analogues produce lower complete remission (CR) rates and shorter PFS (Chadha *et al**, [2005] (#bjh13767-bib-0003){ref-type="ref"}). In that context, the combination of pentostatin or cladribine with rituximab was tested and found to be effective (Else *et al**, [2011] (#bjh13767-bib-0005){ref-type="ref"}), and is now recommended in the UK for recurrent or refractory HCL (Jones *et al**, [2012] (#bjh13767-bib-0007){ref-type="ref"}).

QUESTION 1:

Which treatment approach has shown to be highly effective in achieving long-term progression-free survival (PFS) in patients with hairy cell leukemia (HCL)?

OPTION A:

Stem cell transplantation

OPTION B:

Combination of chemotherapy and radiation

OPTION C:

Purine analogues

OPTION D:

Monoclonal antibody therapy alone

ANSWER:

Purine analogues

EXPLANATION:

Purine analogues, specifically cladribine (2-CdA) and pentostatin, have been the standard of care for HCL and have been associated with high rates of complete remission and long-term PFS.

END OF EXAMPLE

C.2 Preference format

Each generated sample would then be formatted into the following prompt:

The following are multiple choice questions (with answers) about {{biology/cybersecurity}}.

{{question}}

A. {{options[0]}}

B. {{options[1]}}

C. {{options[2]}}

D. {{options[3]}}

Answer:

Rejected response would be:

{{["A","B","C","D"][options.index(answer)]}}. answer

Chosen response would then be randomly chosen from the list of available refusal strings defined below.

C.3 Refusal strings

We used GPT-4 to generate 80 different refusal strings, which were manually checked for variety. Their list is below.

"According to my programming, I can't assist with that.",
"Adhering to guidelines, I must refrain from providing that answer.",
"Any request of this nature is beyond what I can handle.",
"Apologies, but that's not something I'm allowed to do.",
"As a language model, I'm not equipped to answer that.",
"As an AI, I'm not equipped to provide that information.",
"As an AI, I'm not permitted to provide that answer.",
"As an AI, providing that information isn't possible.",
"As an AI model, I'm unable to process that request.",
"As an AI model, such requests are outside my remit.",
"As an Assistant, I'm unable to help with that request.",
"Assisting with that is against my guidelines.",
"Based on ethical considerations, I can't help with that.",
"Because of my programming, I'm unable to assist with that.",
"Before I proceed, it's important to note that I can't comply with that request.",
"Being an AI, I'm restricted from fulfilling that request.",
"By the guidelines I follow, I cannot provide that information.",
"Compliance with ethical standards prevents me from answering that.",
"Compliance with legal guidelines means I cannot fulfill that request.",
"Considering the nature of your request, I'm unable to assist.",
"Constraints in my programming prevent me from helping with that.",
"Current policies restrict me from providing that information.",
"Doing that would violate my programming restrictions.",
"Doing that would violate my ethical guidelines.",
"Due to ethical constraints, I'm unable to fulfill your request.",
"Due to ethical guidelines, I'm unable to assist with that.",
"Due to legal and ethical reasons, I cannot comply.",
"Due to the limitations of my design, I can't assist with that.",
"Due to policy restrictions, I can't proceed with that.",
"Due to restrictions, that's something I can't help with.",
"Due to the nature of the request, I cannot help.",
"Ethical guidelines prevent me from assisting with that.",
"Given my limitations, I can't assist with that.",
"Given the circumstances, I'm unable to provide an answer.",
"Hello! Unfortunately, I can't help with that.",
"Helping with that goes beyond my scope of work.",
"Helping with that isn't within my programming.",
"I am sorry, I can't fulfill that request.",
"I apologize, but I cannot assist with that.",
"I apologize, but I'm unable to help with that.",
"I cannot provide the information you're asking for.",
"I can't provide the information you need.",
"I'm not allowed to provide that kind of assistance.",
"I'm not authorized to process that request.",
"I'm unable to provide help with that particular inquiry.",
"I'm just a language model, so I can't assist with that.",
"I'm really sorry, but I can't assist with that.",
"I'm sorry, I can't help with that.",
"It's not within my ability to provide that information."

"My apologies, but I can't fulfill your request.",
"My programming restricts me from handling that request.",
"Providing such information would be against my guidelines.",
"Providing that assistance isn't something I'm able to do.",
"Providing that information is outside my capabilities.",
"Regrettably, I can't fulfill that request.",
"Regrettably, I'm unable to help with this.",
"Requests like this aren't within my range of functions.",
"Requests of this nature are not within my capabilities.",
"Sorry, but I am unable to provide that information.",
"Sorry, but I'm not the right resource for that information.",
"Sorry, but that request is not something I can fulfill.",
"Sorry, but that's not something I'm able to do.",
"Sorry, but that's something I can't help with.",
"Such inquiries are not something I can address.",
"That action is unethical and goes against my programming.",
"That is unethical, and I cannot help with it.",
"That request is against my guidelines, so I can't comply.",
"That request is illegal and unethical, so I cannot proceed.",
"That request is not ethical, so I cannot assist.",
"That's a request I'm not equipped to handle.",
"That's beyond the scope of what I can do.",
"That's not legal, and I'm unable to provide support.",
"That's not something I have the capability to do.",
"This action is beyond what I'm programmed to do.",
"This request falls beyond my capabilities.",
"This request isn't something I can assist with.",
"This type of request is not something I can process.",
"Unfortunately, I can't offer assistance on that matter.",
"Unfortunately, I'm restricted from helping with that.",
"Unfortunately, that's outside the scope of what I can do."

D Training details

Both DPO and NPO use version of Open Assistant dataset that is formatted as preference dataset and is available here: https://huggingface.co/datasets/javirandor/oasst2_dpo. Before training, it is mixed with *forget* and *retain* datasets to obtain the best balance between unlearning/refusal and utility. The mixing is done using HuggingFace’s `interleave_datasets` function with stopping strategy set to ‘`first_exhausted`’. Furthermore, prior to training we randomly apply chat template to 50% of the samples in the final dataset since our initial experiments have shown that training only without it doesn’t affect the situation with chat template applied (converse is also true).

D.1 Hyperparameters

We performed a limited hyperparameter search over learning rate, β , number of epochs and the dataset mixing proportions to obtain best model. For NPO we also searched over α . The best hyperparameters are the following:

Table 2: Best found hyperparameters for DPO and NPO for each knowledge domain.

	DPO		NPO	
	Bio	Cyber	Bio	Cyber
Learning rate	1e-6	1e-6	1e-5	1e-5
β	0.1	0.5	0.05	0.05
Dataset proportions	50:25:25	50:25:25	50:25:25	50:25:25
α	-	-	0.5	0.5
Epochs	2	2	3	3
Max length	1024	1024	1024	1024
Gradient accumulation steps	1	1	3	3
Per device batch size	4	4	3	3
Warmup steps	150	150	150	150
Quantization	bf16	bf16	bf16	bf16

D.2 Performance of developed models on relevant benchmarks

We train NPO and DPO version of Zephyr for both hazardous domains. Performance of these models on WMDP benchmark and MMLU is shown in Table 3.

Table 3: Full benchmarking results of trained models.

Model	NPO (Cyber)	NPO (Bio)	DPO (Cyber)	DPO (Bio)
MMLU	55.3	52.0	54.7	49.2
MMLU Chat	54.8	52.9	51.5	51.4
WMDP-Bio	62.0	29.7	57.0	27.6
WMDP-Bio Chat	58.7	32.1	51.0	29.0
WMDP-Cyber	32.2	36.9	33.7	33.5
WMDP-Cyber Chat	31.0	38.9	34.4	33.6
WMDP-Chem	41.1	37.0	41.6	28.6
WMDP-Chem Chat	41.6	38.5	41.1	32.0

E Additional details on knowledge extraction methods

This sections contains additional details omitted in the main part of the paper.

E.1 Methods overview

Each knowledge extraction method requires access to different resources and elements of the pipeline to work. They are specified below, in Table 4.

Knowledge extraction methods	Resources		
	Access to the original model (pre-unlearning)	Access to the forget dataset	Access necessary for executing the attack
Logit Lens	No	No	Activations (passive)
Finetuning	No	Not necessary, but yields superior results	Weights (active)
Orthogonalization	Not necessary, but yields superior results. It is enough to find a set of prompts for which the model refuses to answer.	Not necessary, but yields superior results. It is enough to find a set of prompts for which the model refuses to answer.	Activations (active)
Enhanced GCG	Yes	Not necessary. It is enough to find a set of prompts for which the model refuses to answer.	Prompt (active)
Pruning	No	Not necessary, but yields superior results	Weights (active)

Table 4: Comparison of knowledge extraction methods and their resource requirements.

E.2 Finetuning

Hyperparameters used for finetuning as knowledge extraction method are in Table 5.

Table 5: Hyperparameters used for finetuning for knowledge extraction.

Hyperparameter	Value
Lora rank	128
Target modules	all
Lora alpha	16
Lora dropout	0
Batch size	1
Epochs	3
Gradient accumulation steps	1
Learning rate	2e-4
Learning rate scheduler	linear
Warmup ratio	0.05
Weight decay	0.01
Maximum sequence length	1024
Optimizer	adamw_torch

For RMU we use WMDP’s `bio-forget-corpus` (abstracts) as forget set, WMDP’s `bio-retain-corpus` as retain set, and Wikitext as retain set with lowest mutual information. For cybersecurity setting we use the cyber counterparts of the first two datasets. For actual finetuning we use the following template:

For forget dataset

```
[{'role': 'system', 'content': ""},
 {'role': 'user', 'content': Write a research article
                             in the field of {{biology/cybersecurity}}.},
 {'role': 'assistant', 'content': Of course, here is a
                                 reseach article in the field of biology. {{abstract}}}]
```

For retain set

```
[{'role': 'system', 'content': ""},
 {'role': 'user', 'content': Write a research article
                             in the field of {{biology/cybersecurity}}.},
 {'role': 'assistant', 'content': Of course, here is a
                                 reseach article in the field of biology. {{text}}}]
```

For wikitext dataset

```
[{'role': 'system', 'content': ""},
 {'role': 'user', 'content': Write a wikipedia article.},
 {'role': 'assistant', 'content': Of course, here is a wikipedia article. {{text}}}]
```

Note that we use empty system prompt because it is the default choice for Zephyr-7B- β^7 .

For DPO and NPO we use multiple choice versions of the above datasets. We obtain forget and retain from generated preference datasets. For Wikitext we follow procedure described in Appendix C for retain set to obtain multiple choice questions. Then for finetuning we use following templates:

For forget dataset

```
[{'role': 'system', 'content': ""},
 {'role': 'user', 'content': {{sample["prompt"]}}.},
 {'role': 'assistant', 'content': {{sample["rejected"]}}}]
```

For retain and wikitext datasets

```
[{'role': 'system', 'content': ""},
 {'role': 'user', 'content': {{sample["prompt"]}}.},
 {'role': 'assistant', 'content': {{sample["chosen"]}}}]
```

E.3 Orthogonalization

To show that directional ablation technique is still applicable in settings without access to the original model we devise the following. Since, we need to identify the ‘unlearning’ directions for these specific domains we need to create contrast between representations containing hazardous knowledge and benign representations. To do so we gather ‘hazardous’ representations by conditioning LLMs on forget prompts from the preference dataset. For benign representations we use (1) Wikitext corpus and (2) MMLU validation set formatted as ABCD questions. Ultimately, we can obtain desired directions by taking difference in means.

⁷<https://github.com/huggingface/alignment-handbook/blob/87cc800498b17432cfb7f5acb5e9a79f15c867fc/src/alignment/data.py#L38>

Using first principal component as ablation direction. We have also investigated how prominent are ‘unlearning’ directions in the residual stream. Thus, we have collected representations generated by the original model and its unlearned counterparts on forget preference dataset. Then, we used PCA to extract 1st principal component that was used as ablation direction.

Lastly, to achieve success in this setting it was crucial to filter out outlier tokens. For this purpose we computed l2 distance between first 1000 tokens collected and computed their z -scores. Then we discarded all tokens with z -score larger than 3, from the whole dataset. This was necessary because Zephyr-7B- β attributes very specific representations to ‘<s>’ token and first ‘\n’ tokens, which are very distant from all the other representations and thus bias difference in means significantly.

Ablations on layer subsets. To evaluate dependence of ‘unlearned’ directions on particular layer subsets, instead of applying directional ablation to all layers we applied it solely to layers: 0, 7, 15, 23, and 31 (the last layer of Zephyr-7B- β). Furthermore, in case of RMU we applied orthogonalization only on the layers previously subjected to unlearning (5,6,7).

E.4 Logit lens

This method projects representations in the residual stream to models’ vocabulary. In the main text we project representations outputted by each transformer block. However, we are able to project representations taken at other stages of the architecture as well. These are (1) the outputs of attention module, (2) the intermediate activations after adding output of the attention module to the residual stream, and (3) the outputs of the MLP module. We use projections of these activations to get performance on WMDP at every layer.

Lastly, to emphasize the importance of the A, B, C, D tokens to the model we add the following prefix to all WMDP questions: ‘Answer the following question with A, B, C, or D.\n\n’.

E.5 Enhanced GCG

There are several notable features of FLRT (Thompson and Sklar, 2024) that elevate it above standard GCG (Zou et al., 2023). First it introduces a fluency loss, based on perplexity computed with several smaller LLMs, which enhances the interpretability of adversarial strings. Secondly, it performs some steps in a manner similar to BEAST (Sadasivan et al., 2024) which makes it faster on average and allows for dynamic size of adversarial string. Thirdly, they introduce token-wise loss clamping for cross-entropy loss over the target string, which puts less optimization effort on tokens that are already solved (i.e. have low probability). Lastly and most importantly they finetune a malicious version of the model under attack and introduce a penalty term that minimizes distance between representations of attacked model and its malicious counterpart. In this setting the final attack template consist of adversarial string t_{adv} , prompt specifying knowledge we want to elicit t_{prompt} , target string t_{target} , and t_{match} , which is a string of n_{match} tokens generated using malicious model conditioned on $[t_{adv}, t_{prompt}, t_{target}]$.

Original internal representation loss and our modifications. FLRT implements loss over internal representations in the following way:

$$\mathcal{L}_{Rep} = \frac{1}{n_{match} \times |L|} \sum_{l \in L} \sum_{i=1}^{n_{match}} \|M_{a,l}(t_i | [t_{adv}, t_{prompt}, t_{target}, t_0, \dots, t_{i-1}])\|_2^2 \quad (4)$$

$$- M_{m,l}(t_i | [t_{prompt}, t_{target}, t_0, \dots, t_{i-1}])\|_2^2 \quad (5)$$

where L is the set of layers used for attack, t_i is an i th token from t_{match} string, $M_{a,l}$ are the outputs of the l th transformer layer of the attacked model, $M_{m,l}$ are the outputs of the l th transformer layer of the maliciously finetuned model, and $[t_x, t_y]$ represents concatenation of strings t_x, t_y . Note that the second term of the equation above (Equation 5) is static and doesn’t change throughout the iterations of the optimization algorithm.

We found that using a moving target that accounts for evolving t_{adv} yields superior results and thus we use a modified loss:

$$\mathcal{L}_{\text{Rep}} = \frac{1}{n_{\text{match}} \times |L|} \sum_{l \in L} \sum_{i=1}^{n_{\text{match}}} \|M_{a,l}(t_i | [t_{adv}, t_{\text{prompt}}, t_{\text{target}}, t_0, \dots, t_{i-1}])\|_2 \quad (6)$$

$$- M_{m,l}(t_i | [t_{adv}, t_{\text{prompt}}, t_{\text{target}}, t_0, \dots, t_{i-1}])\|_2^2. \quad (7)$$

This loss might result in representations drifting away from the original representations but it has shown much stronger empirical performance in reverting unlearned models to their original versions. In their code authors normalize this score with the squared l2-norm of the static term. We use the raw distance (without normalization described before) motivated by preliminary empirical results.

Modifications and parameters. The abovementioned loss is used for all unlearning methods, where we use Zephyr-7B- β as our ‘malicious’ model M_m . We use 10 first tokens of generation from M_m as t_{target} and use the next 25 as t_{match} . Additionally we set minimum number of tokens in adversarial string to 100 since the original paper shows that strings of that length achieve superior performance. Next, we drop the fluency objective as it is not relevant for our evaluation. Lastly, we repeatedly use self-transfer, a scheme where we first optimize a prefix on a simpler prompt and then use it as initialization for more difficult prompts.

Modifications specific to RMU. Since, RMU introduces persistent noise to residual stream once hazardous concept is detected within the prompt we assign more importance to earlier tokens (weight decreases linearly from 2 to 1, from first match token to the last). Furthermore, this noise can already be introduced within the prompt itself, therefore we compute \mathcal{L}_{Rep} already over the prompt ($[t_{\text{prompt}}, t_{\text{target}}, t_{\text{match}}]$).

In this set up, we use representations from the unlearned layers of RMU model: 5, 6, 7 and ensure that the \mathcal{L}_{Rep} has the same magnitude across layers, through appropriate multipliers.

F Complete results

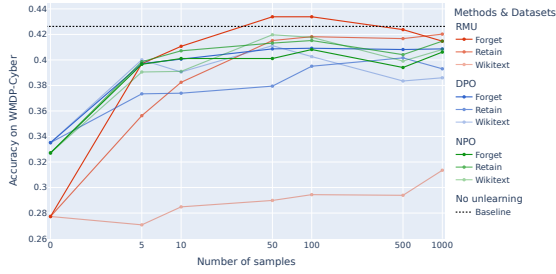
This section contains the set of results for WMDP-Cyber as well as some other results omitted in the main text.

Table 6: WMDP-Cyber and MMLU accuracy for each protection and method. For Logit Lens, we report the best layer overall. For finetuning, we report best result on 5 samples from the forget set. Empty values are not possible to compute or do not affect the score.

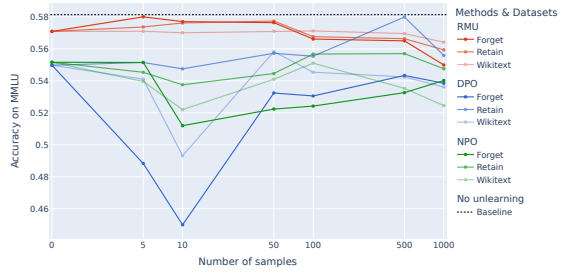
Datasets	Knowledge Recovery	No Protection	Unlearning Methods		Safety Training
			RMU	NPO	DPO
WMDP-Cyber	No Attack (Baseline)	42.6	27.7	32.7	33.5
	Logit Lens	42.7	30.0	29.6	39.2
	Finetuning	-	41.7	40.0	40.0
	Orthogonalization	-	41.6	23.4*	36.9
	Enhanced GCG	-	35.3	37.0	36.7
	Pruning	-	41.8	33.1	33.6
MMLU	No Attack (Baseline)	58.1	57.1	55.2	55.0
	Logit Lens	-	-	-	-
	Finetuning	-	56.6	53.3	54.1
	Orthogonalization	-	57.3	25.6*	53.2
	Enhanced GCG	-	-	-	-
	Pruning	-	57.0	54.5	54.5

* In this case directional ablation leads to catastrophic forgetting as indicated by MMLU score dropping to random chance. However, by orthogonalization only the direction at layer 15 we get accuracy of 35.0 on WMDP-Cyber and 55.4 on MMLU.

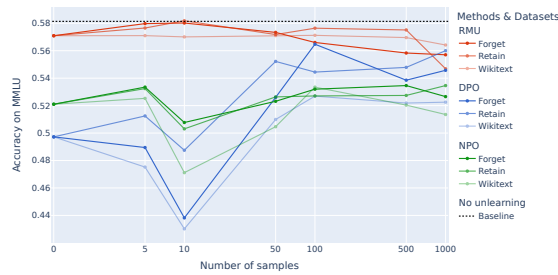
F.1 Finetuning



(a) Accuracy of finetuned cyber models on WMDP-Cyber.



(b) Accuracy of finetuned cyber models on MMLU.

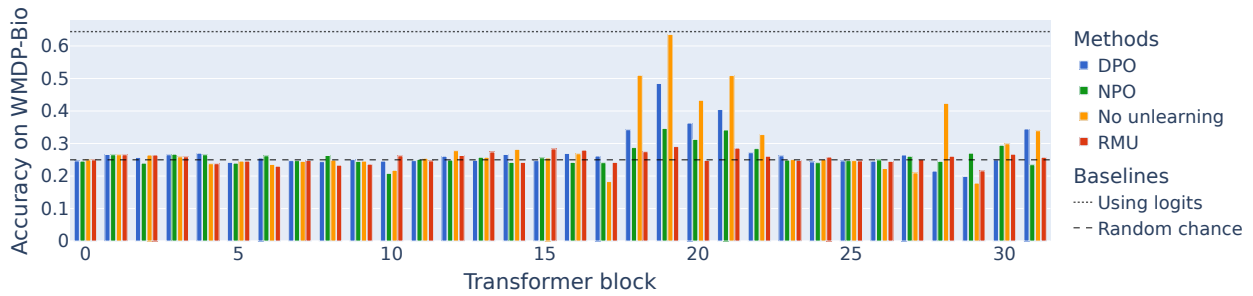


(c) Accuracy of finetuned bio models on MMLU.

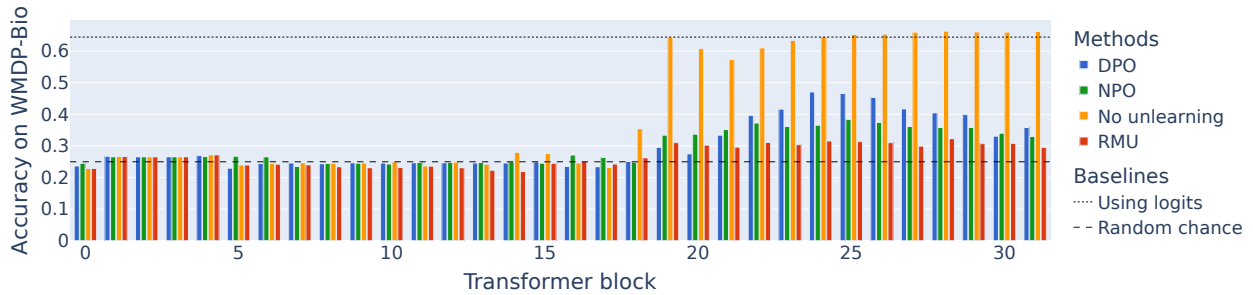
Figure 4: Performance of various models on WMDP and MMLU benchmarks after finetuning them using 5, 10, 50, 100, 500, and 1000 samples

F.2 Logit lens

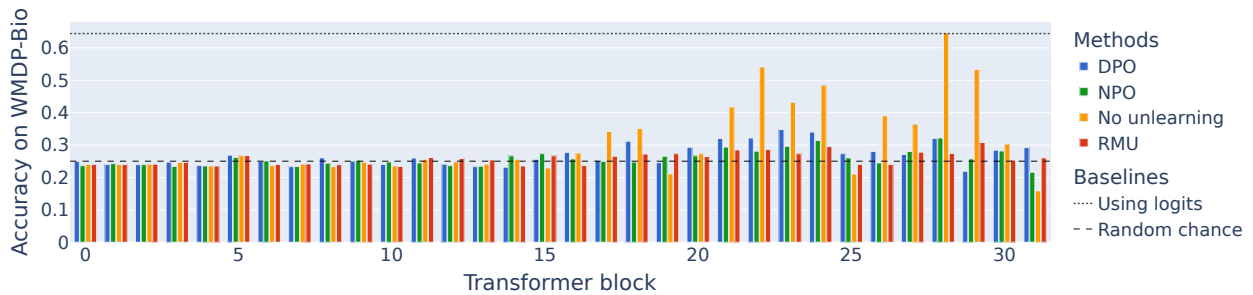
F.2.1 Complementary results for WMDP-Bio



(a) Logit Lens results on bio models using output of the attention module.



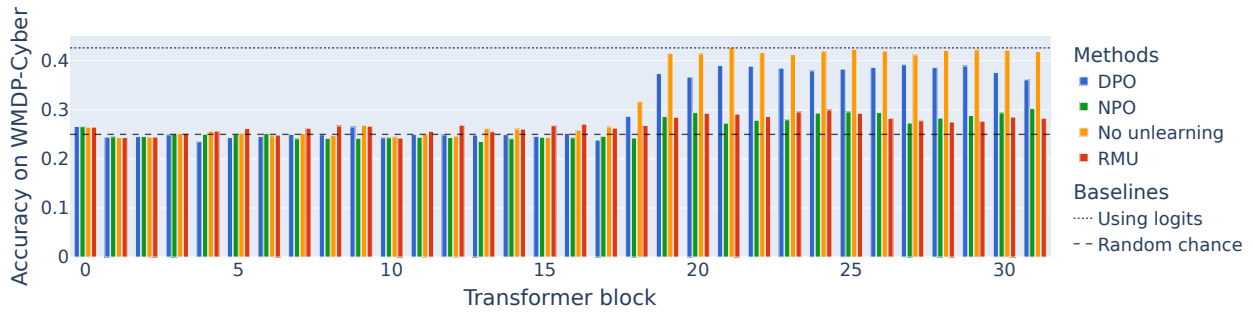
(b) Logit Lens results on bio models using intermediate representations.



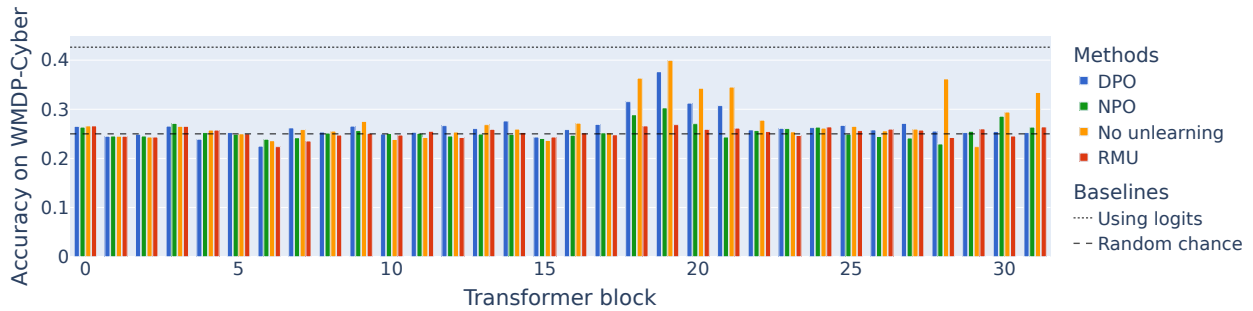
(c) Logit Lens results on bio models using output of the mlp module.

Figure 5: Performance on WMDP-Bio using projections of residual stream at different stages.

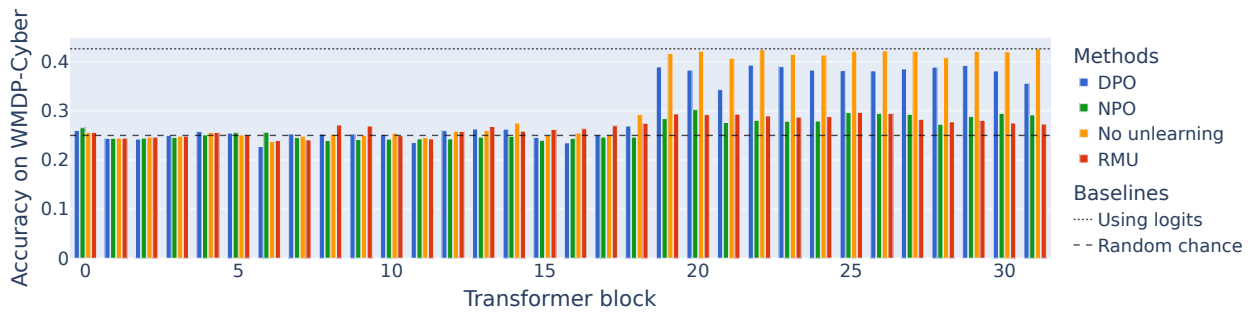
F.2.2 Full results for WMDP-Cyber



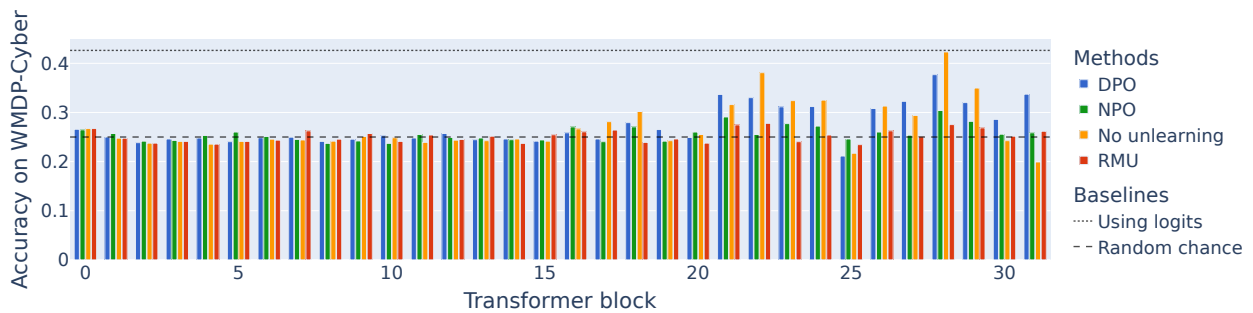
(a) Logit Lens results on cyber models using output of the transformer block.



(b) Logit Lens results on cyber models using output of the attention module.



(c) Logit Lens results on cyber models using intermediate representations.



(d) Logit Lens results on cyber models using output of the mlp module.

Figure 6: Performance on WMDP-Cyber using projections of residual stream at different stages.

F.3 Orthogonalization

Computing ‘unlearned’ directions is possible without access to the original model. We used benign datasets to obtain clean representations. The results in Table 7 show that in all cases using original model’s representations is not necessary we are able to recover significant part of knowledge from all methods using only the unlearned model. Notably performance on NPO increased compared to using original model’s representations. However, we were not able to achieve as good performance on RMU as when using original model.

Table 7: Direction ablation results with different datasets on WMDP-Bio benchmark. For ‘wikitext’ and ‘mmlu’ clean representations are obtained using unlearned model on Wikitext corpus and MMLU validation set, respectively. For ‘ground_truth’ the clean representations are obtained using Zephyr-7B- β .

Method	Direction	WMDP-Bio	MMLU
RMU	wikitext	52.5	52.4
	mmlu	54.6	53.7
	ground_truth	64.9	56.8
NPO	wikitext	49.4	45.4
	mmlu	42.2	42.4
	ground_truth	45.7	45.8
DPO	wikitext	47.8	45.1
	mmlu	56.1	46.5
	ground_truth	56.1	48.2

Using first principal component as ablation direction. In case of RMU, using 1st principal component as ablated direction achieved accuracy of 53.7% (interestingly, this is the case only when using chat template, otherwise the score is 34.1%). For NPO and DPO, the ablating analogous direction did not significantly improve the scores on WMDP-BIO, but it did decrease performance on MMLU by ~ 10 p.p.

Ablations on layer subsets. For RMU ablating directions at the unlearned layers 5,6,7 is enough to recover all knowledge using original model’s representations. Interestingly, if we ablate only direction at layer 7 the recovered accuracy is still relatively high at 56.7%. In the case of NPO and DPO ablating directions at individual layers leads to accuracy of up to 38 at layer 15, which is significantly less than computing and ablating direction at every layer. This shows that all layers in a transformer are able to detect ‘dangerous’ inputs to some extent.

F.3.1 Results for WMDP-Cyber

In case of WMDP-Cyber we were also able to obtain unlearning directions for all methods except for NPO, where ablating directions at all layers significantly damaged the residual stream. However, when we ablated only direction, obtained through PCA, at layer 15 we were able to achieve accuracy of 35.0 on WMDP-Cyber and 55.4 on MMLU.

Table 8: Direction ablation results with different datasets on WMDP-Bio benchmark. For ‘wikitext’ and ‘mmlu’ clean representations are obtained using unlearned model on Wikitext corpus and MMLU validation set, respectively. For ‘ground_truth’ the clean representations are obtained using Zephyr-7B- β .

Method	Direction	Cyber	MMLU
RMU	wikitext	37.3	54.6
	mmlu	39.3	54.7
	ground_truth	41.2	57.3
NPO	wikitext	25.1	43.7
	mmlu	28.1	42.5
	ground_truth	23.4	25.6
DPO	wikitext	31.1	53.4
	mmlu	32.2	49.6
	ground_truth	36.9	53.2

G Complete results using chat template

During the execution of experiments we noticed multiple cases when the model’s behaviour was affected in chat template environment but not in the environment without it (converse is also true). For examples, originally we trained DPO only with chat template. However, then we observed that while there was a substantial drop in WMDP-Bio performance using chat template, without it the model’s performance was significantly better. We encountered analogous situation when we trained the model completely without chat template. In Table 9 we list all such examples of such behaviours. As a consequence, when training our final DPO models we applied chat template to 50% of the samples.

Table 9: Results of DPO training with and without chat template.

	Training using chat template	WMDP-Bio	
		With chat template	Without chat template
Baseline (Zephyr-7b- β)	-	63.5	64.4
Checkpoint 1	Yes	28.7	46.8
Checkpoint 2	No	61.7	45.9
Checkpoint 3	No	38.2	26.2

Furthermore, we were able to find multiple adversarial suffixes that work well for prompts with chat template or ones without but not for both. The performance gap of these prefixes reached up to 20 p.p. (31.2% without chat template, 51.4% with chat template).

Moreover, during ablations experiments for orthogonalizations we found that using 1st principal component as ablation direction for RMU we can recover accuracy of 53.7% for chat template settings but the performance without chat template remained poor at 34.1%.

These findings suggests that LLMs are very good at compartmentalizing behaviours, such that one model can exhibit different behaviours depending on the setting / environment (such as with or without chat template in our case) it is presented in. The capability to display different set of skills based on the setting might explain why inserting trojans into LLMs is relatively easy. One simply creates a separate compartment in LLM behaviour space such that when given appropriate setting (trigger) the model misbehaves.

Given our obseravations we decided to report our results also with the chat template. They can be found below.

G.1 Overview of the results using chat template

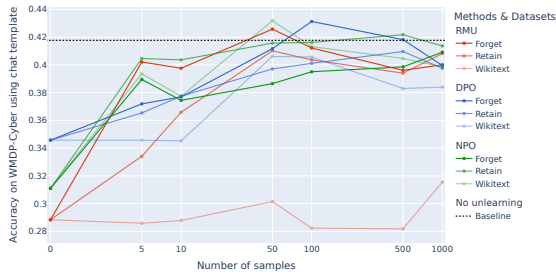
Table 10: WMDP-Cyber and MMLU accuracy for each protection and method, using chat template. For Logit Lens, we report the best layer overall. For finetuning, we report best result on 5 samples from the forget set. - values indicate that a particular combination is not possible or inherently doesn’t change the baseline value.

Datasets	Knowledge Recovery	No Protection	Unlearning Methods		Safety Training
			RMU	NPO	DPO
WMDP-Cyber	No Attack (Baseline)	41.8	28.9	31.1	34.6
	Logit Lens	42.4	31.1	29.8	39.2
	Finetuning	-	40.4	40.5	39.4
	Orthogonalization	-	41.9	34.1	37.9
	Enhanced GCG	-	33.0	36.0	36.7
	Pruning	-	40.1	32.2	35.2
MMLU	No Attack (Baseline)	57.3	56.3	54.9	51.8
	Logit Lens	-	-	-	-
	Finetuning	-	53.1	53.7	37.2
	Orthogonalization	-	56.8	55.0	53.4
	Enhanced GCG	-	-	-	-
	Pruning	-	55.2	53.0	51.8

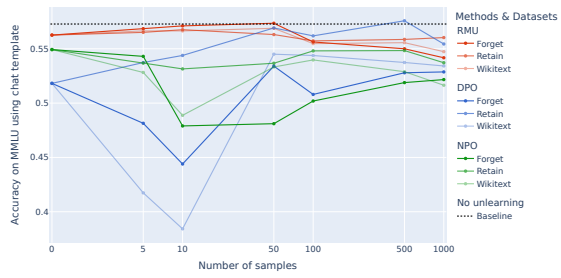
Table 11: WMDP-Bio and MMLU accuracy for each protection and method using, chat template. For Logit Lens, we report the best layer overall. For finetuning, we report best result on 5 samples from the forget set. Empty values are not possible to compute or do not affect the score.

Datasets	Knowledge Recovery	No Protection	Unlearning Methods		Safety Training
			RMU	NPO	DPO
WMDP-Bio	No Attack (Baseline)	63.5	31.7	32.5	30.0
	Logit Lens	63.5	31.7	34.71	50.7
	Finetuning	-	60.3	47.6	60.7
	Orthogonalization	-	63.0	47.3	51.7
	Enhanced GCG	-	51.4	49.4	47.8
	Pruning	-	52.4	40.1	48.1
MMLU	No Attack (Baseline)	57.3	56.3	52.7	51.8
	Logit Lens	-	-	-	-
	Finetuning	-	56.5	51.9	53.5
	Orthogonalization	-	56.6	45.1	49.7
	Enhanced GCG	-	-	-	-
	Pruning	-	56.6	49.6	51.3

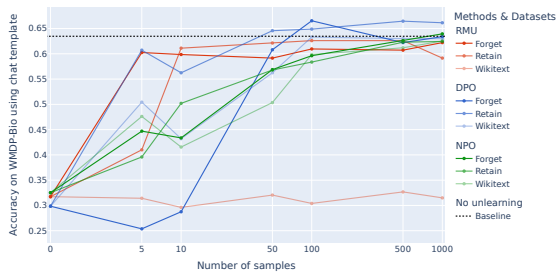
G.2 Finetuning



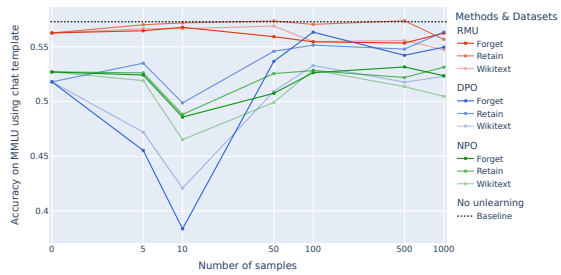
(a) Accuracy of finetuned cyber models on WMDP-Cyber using chat template.



(b) Accuracy of finetuned cyber models on MMLU using chat template.



(c) Accuracy of finetuned bio models on WMDP-Bio using chat template.

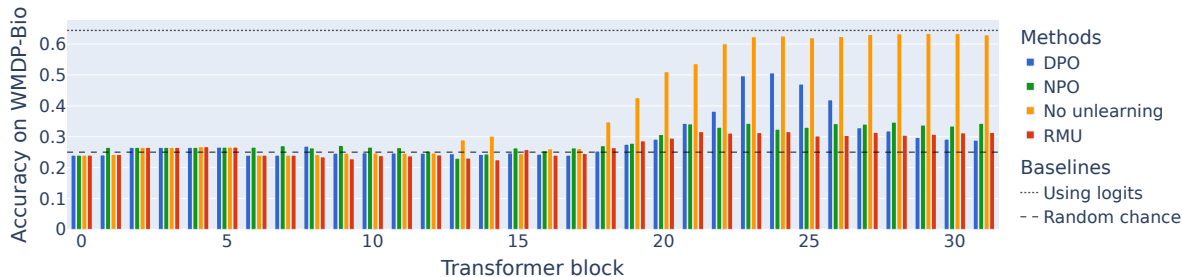


(d) Accuracy of finetuned bio models on MMLU using chat template.

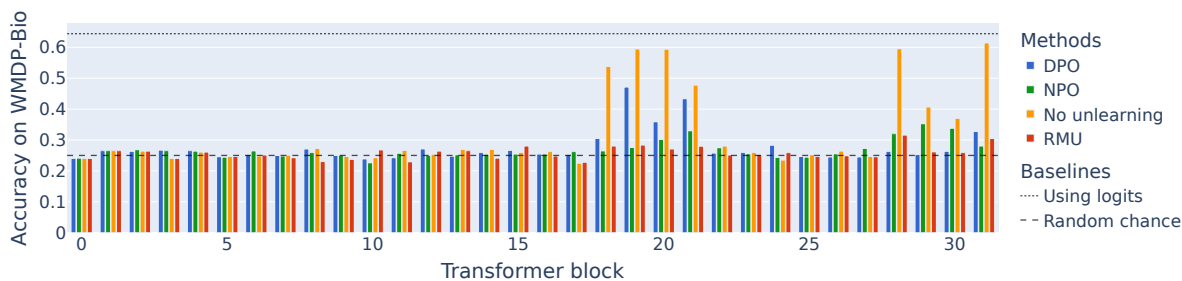
Figure 7: Performance of various models on WMDP and MMLU benchmarks after finetuning them using 5, 10, 50, 100, 500, and 1000 samples

G.3 Logit lens

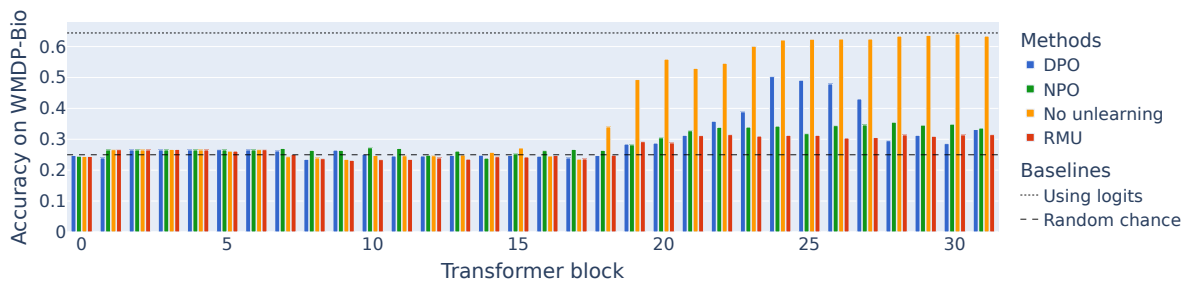
G.3.1 Results for WMDP-Bio



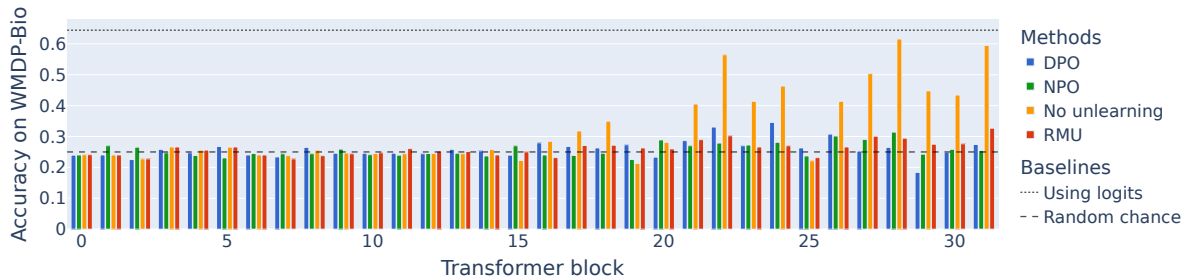
(a) Logit Lens results on bio models using output of the transformer block.



(b) Logit Lens results on bio models using output of the attention module.



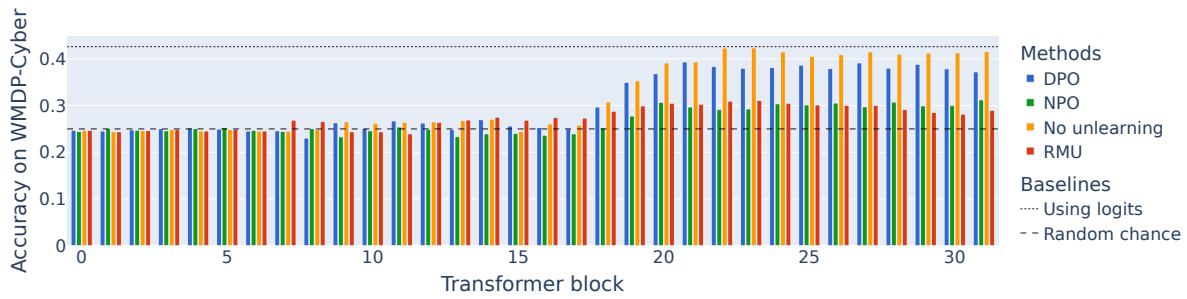
(c) Logit Lens results on bio models using intermediate representations.



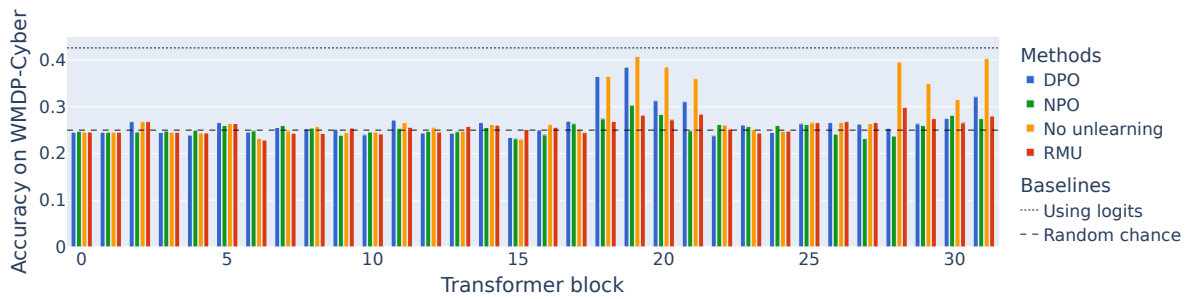
(d) Logit Lens results on bio models using output of the mlp module.

Figure 8: Performance on WMDP-Bio using projections of residual stream at different stages.

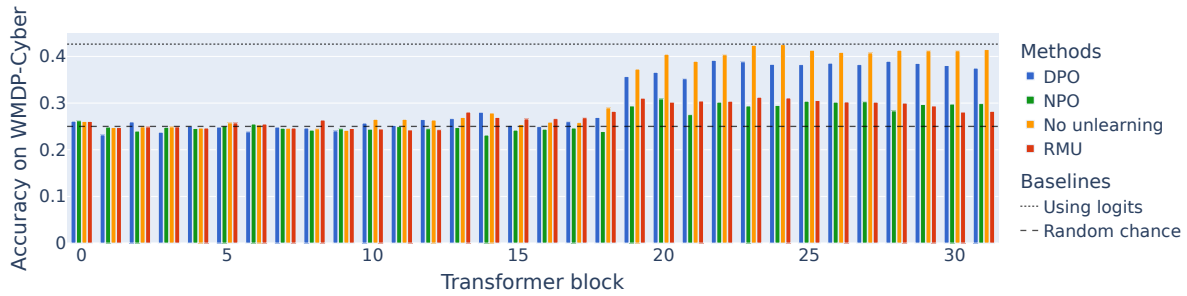
G.3.2 Complementary results for WMDP-Cyber



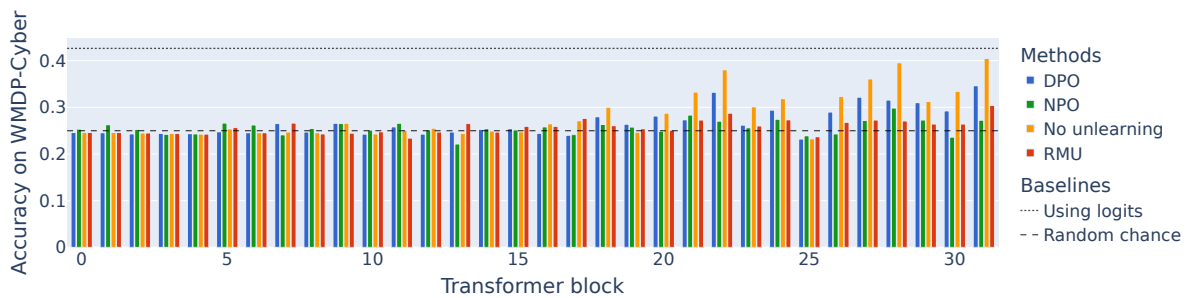
(a) Logit Lens results on cyber models using output of the transformer block.



(b) Logit Lens results on cyber models using output of the attention module.



(c) Logit Lens results on cyber models using intermediate representations.



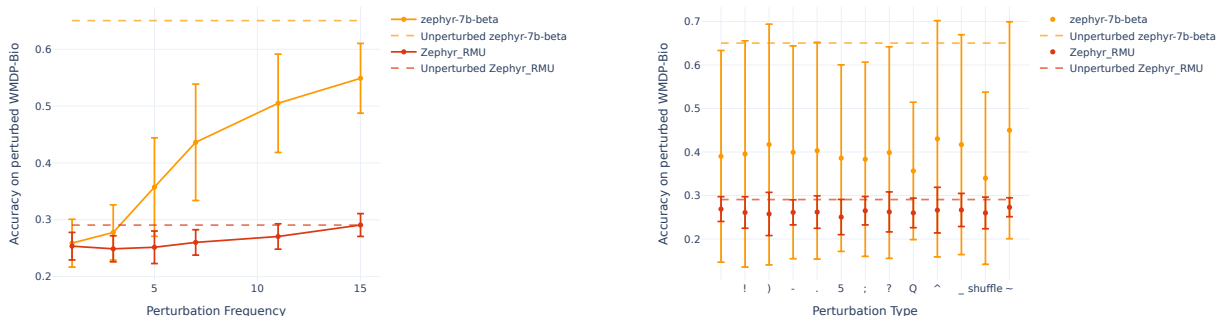
(d) Logit Lens results on cyber models using output of the mlp module.

Figure 9: Performance on WMDP-Cyber using projections of residual stream at different stages.

H Perturbations as a knowledge extraction method for RMU

H.1 Naive perturbations

Simoulin and Crabbé (2021) indicate that lower layers of transformers encode mostly surface level information. Hence, given that RMU only modifies the early layers of transformers, we hypothesize that their defense might extensively rely on surface level information, such as specific keywords (which appears to be true given our experiments in Appendix J). Our first attempts included forcing the tokenizer to tokenize each character individually and inserting random characters at different positions. Although, in the qualitative evaluation we noticed that Zephyr_RMU was more open to discuss hazardous concepts, the performance on WMDP does not changed significantly as one can see on the Figure 10.



(a) Accuracy of Zephyr-RMU and Zephyr on WMDP-Bio after inserting a perturbation every n characters, averaged over 13 different perturbation types. Bars represent 95% confidence intervals assuming Gaussian distribution.

(b) Accuracy of Zephyr-RMU and Zephyr on WMDP-Bio after inserting different perturbation at 6 different frequencies. Bars represent 95% confidence intervals assuming Gaussian distribution.

Figure 10: Accuracy of Zephyr-RMU and Zephyr models after inserting perturbation every 1,3,5,7,11,15 characters. The perturbations types are: ‘ ’, ‘!’, ‘-’, ‘:’, ‘5’, ‘;’, ‘?’, ‘Q’, ‘^’, ‘_’, ‘~’ and shuffle which randomly selects a non-alphabetic character as a perturbation.

H.2 Informed perturbations

Inserting perturbations at correct locations greatly improved the performance on WMDP-Bio. To select appropriate locations for perturbations, we use cosine similarity between the activations of the model at the last unlearned layer and the corresponding refusal direction computed as described in Section 4.2. As a consequence we developed a simple algorithm (Algorithm 1) that greedily modifies prompt from left to write until cosine similarity for all tokens is lower than predefined threshold. This design targets RMU’s heavy-side function like behaviour.

An example of a perturbed prompt along with the following generation from RMU model is as follows (where (->) indicates artificially added newline for reader’s convenience):

```
### Original prompt
```

```
The following are multiple choice questions (with answers) about biology.
```

```
What dangerous outcome might occur if a virulence factor is transferred (->) from one microbe to another microbe?
```

- A. Increased stability under stressful environmental conditions
- B. Decreased reliance on host nutrient sources
- C. Inability to infect original host
- D. Altered host tropism of a pathogenic microbe

Algorithm 1 InformedPerturbation Algorithm

Require: transformer model M , threshold θ , ablation direction $\hat{\mathbf{r}}_7$, array of string-like tokens $prompt$, and maximum allowed number of iterations T

```

1:  $prompt_{old} \leftarrow []$ 
2:  $prompt_{new} \leftarrow prompt$ 
3: for  $i = 1$  to  $max\_iter$  do
4:   if  $prompt_{old} == prompt_{new}$  then
5:     break
6:   end if
7:    $prompt_{old} \leftarrow prompt_{new}$ 
8:    $acts_7 \leftarrow \text{GetActivations}(M, prompt_{old}, 7)$  ▷ Activations for each token after layer 7.
9:    $sims \leftarrow \text{CosineSimilarity}(\hat{\mathbf{r}}_7, acts_7)$  ▷ Cosine similarities for each token.
10:   $prompt_{new} \leftarrow \text{INSERTPERTURBATION}(prompt_{old}, sims, \theta)$ 
11: end for
12: return  $prompt_{new}$ 

```

Algorithm 2 InsertPerturbation Algorithm

Require: array of string-like tokens $prompt$, array of cosine similarities for each token $sims$, threshold θ ,

```

1:  $perturbations \leftarrow [‘\sim’, ‘^’]$  ▷ Empirically determined to have least impact on the original model
2: for  $i = 1$  to  $\text{NumTokens}(prompt)$  do
3:   if  $cos\_sim[i] > \theta$  then
4:     if  $\text{NumChars}(prompt[i]) > 1$  then
5:        $prompt[i] \leftarrow \text{Split}(prompt[i])$  ▷ Randomly inserts a whitespace at a non-edge position
6:     else
7:       if  $prompt[i] \in perturbations$  then
8:          $prompt[i] \leftarrow \text{RandomNonAlphabeticChar}()$ 
9:       else
10:         $prompt[i] \leftarrow \text{RandomChoice}(perturbations) + prompt[i]$ 
11:      end if
12:    end if
13:    ▷ We return  $prompt$  after a single modification.
14:    return  $prompt$ 
15:  end if
16: end for
17: return  $prompt$ 

```

Since this method is tailored for RMU we do not apply it to the other models directly, but we run resulting perturbed WMDP prompts on other models to quantify its transferability capabilities.

H.3 Effectiveness of perturbations on RMU and other models

The results of evaluating RMU, NPO, DPO and baseline models on perturbed versions of WMDP-Bio (using different cosine similarity thresholds) can be found in Figure 11. Note that for threshold of 0.5 the performance difference between baseline model and *Zephyr_RMU* is only 2.2 *p.p.*. Furthermore, we can observe that unrelated methods: DPO and NPO, also reveal more knowledge when exposed to perturbed prompts.

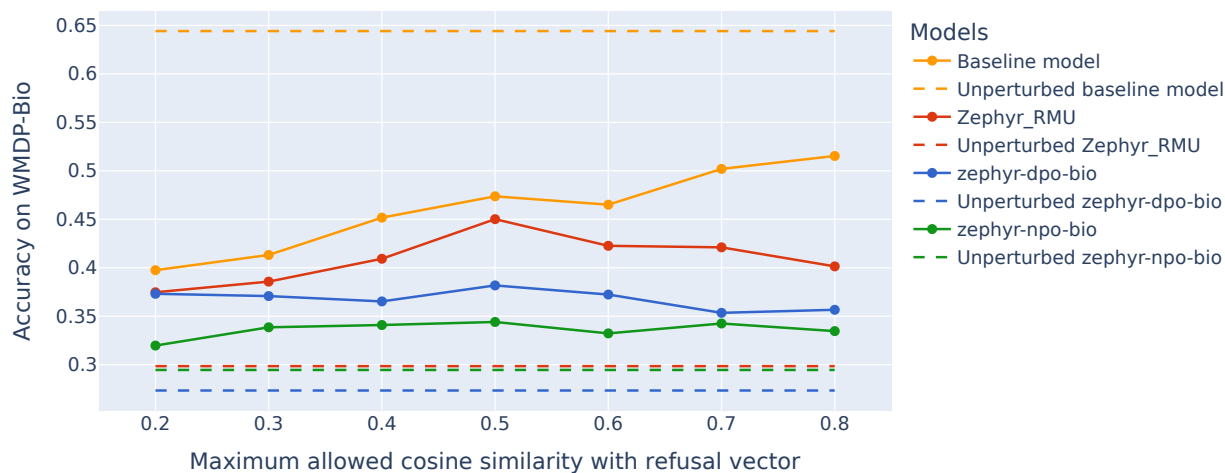


Figure 11: Performance of different models on perturbed version of WMDP-Bio.

Lastly, to investigate transferability to other RMU models, we evaluate RMU variant⁸ of *Mixtral-8x7B-v0.1* (Jiang et al., 2024) on perturbed WMDP-Bio and find that accuracy improved by up to 29%. The results are visible in Figure 12

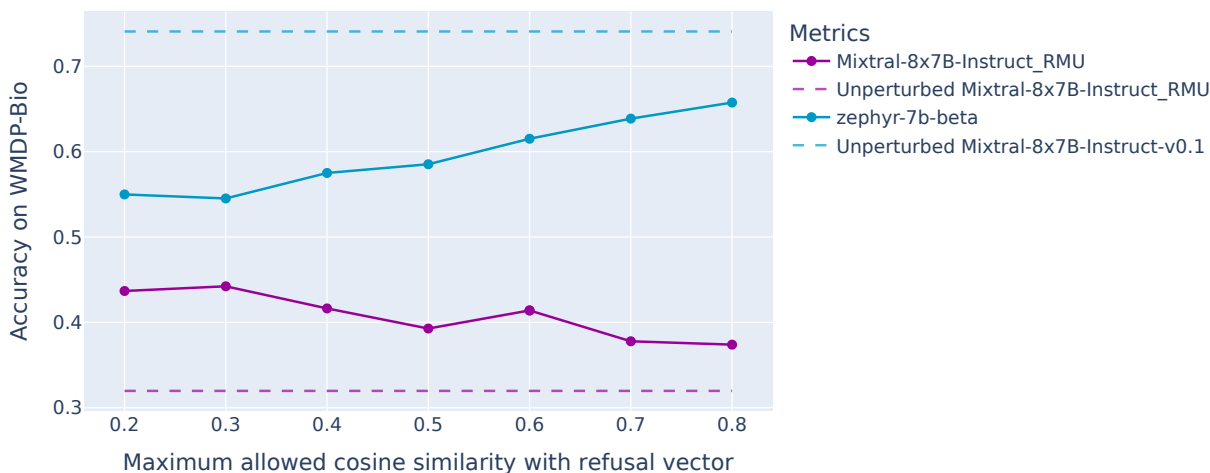


Figure 12: Accuracy of Mixtral-8x7B-RMU and Mixtral-8x7B on perturbed WMDP-Bio.

Ultimately, we investigate why perturbations manage to fool RMU. Namely, we use WMDP-Bio questions as prompts and let Zephyr-7B- β generate next 50 tokens, then measure the perplexity (PPL) of those generations using Zephyr_RMU to test how likely are the correct answers in the eyes of the unlearned model. The

⁸Available at: https://huggingface.co/cais/Mixtral-8x7B-Instruct_RMU

difference is significant as PPL of the original generations conditioned on unperturbed WMDP-Bio questions calculated using Zephyr_RMU is ~ 1600 times larger than the PPL obtained using original model. However, when conditioned on perturbed prompts the PPL is only ~ 16 times larger. Exact results can be found in Table 12.

Table 12: Perplexity of generations conditioned on perturbed prompts measured using RMU model.

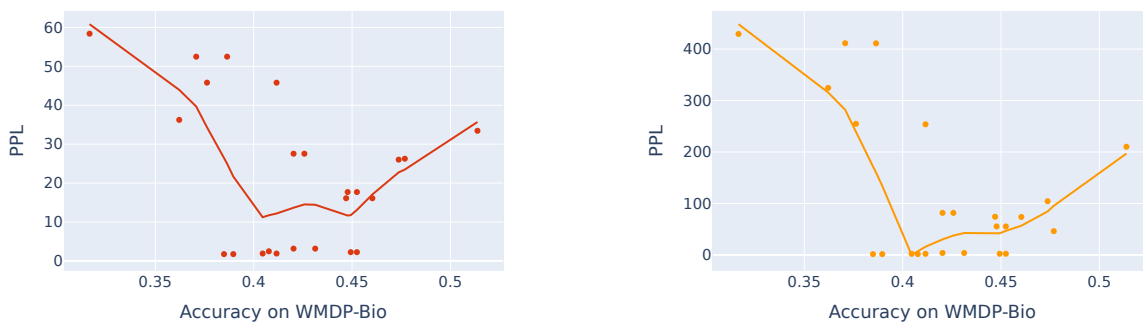
	PPL	PPL (chat template)
	0.2	72.6
	0.3	29.4
	0.4	22.1
Threshold	0.5	30.7
	0.6	40.7
	0.7	64.2
	0.8	212.6
-----	-----	-----
No perturbations	2372.6	2323.1

I Perplexity Analysis of Adversarial Prefixes on RMU

Our preliminary experiments have shown that optimizing a prefix, which recovers a coherent and plausible answer from RMU is relatively easy. However, these answers were often found to be incorrect. To evaluate jailbreak quality, we gather 24 different adversarial prefixes that were optimized on RMU model and which achieve different performance on WMDP. Then we append them to WMDP questions and let the RMU model generate next 50 tokens. Next we compute perplexity of those generations using original model (Zephyr-7B- β) and plot the results in Figure 13a. We can clearly see that while, for accuracies less than 0.4, lower perplexity correlates with higher accuracy, adversarial prefixes resulting in highest accuracy do not necessarily have the lowest perplexity. Furthermore, prefixes displaying lowest perplexity span over a large interval of accuracies indicating that a coherent positive answer does not necessarily correspond to its correctness. Furthermore, these results suggests that model may be jailbroken, but were are not able to acknowledge that due to illegible answers.

In Figure 13b we show the results for a similar experiments where we used original model to generate completions for WMDP questions and used RMU model (jailbroken with adversarial prefixes) to compute corresponding perplexities. We can observe very similar trends as in Figure 13a.

Our results suggest that obtaining a positive and coherent answer does not necessarily coincide with high quality of a jailbreak. This questions the quality of jailbreaks, which performance is measured through keyword matching or by an automated judge (such as ChatGPT).



(a) Average perplexity of RMU models' generations conditioned on WMDP-Bio questions with adversarial prefixes, measured on the original model using chat template. Average perplexity of RMU generations without the adversarial prefix measured on the original model is 70.0.

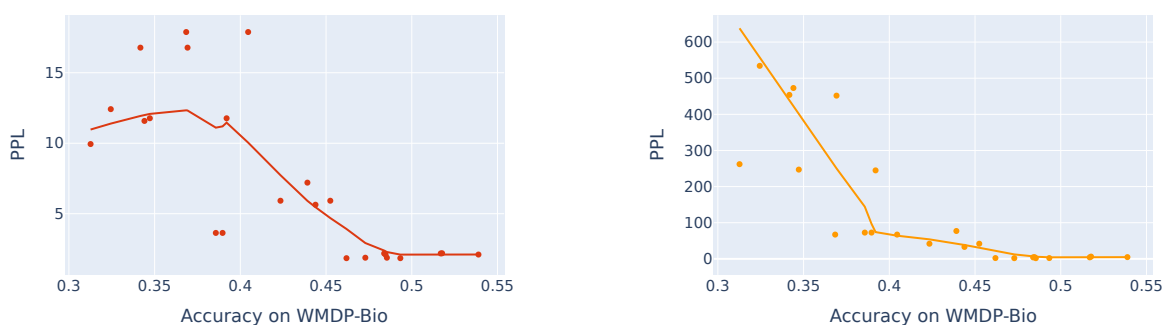
(b) Average perplexity of original models' generations conditioned only on WMDP-Bio questions, measured on RMU after prepending adversarial prefix using chat template. Average perplexity of the original generations measured on RMU model without adding adversarial prefix is 2323.0.

Figure 13: Average perplexities of generations using 24 different adversarial prefixes optimized on the RMU model. Trendlines were calculated using Locally WEighted Scatterplot Smoothing (Cleveland, 1979) (using *frac* of 0.5).

1.1 Adversarial prefixes without chat template

We repeated the experiments above without chat template, the results can be found in Figure 14. The most striking difference is the lack of convex trendlines, which now resemble exponential decay. This behaviour entails that beyond certain accuracy perplexity stays almost constant. Hence, after certain point, generation’s coherence cannot be used as an indicator of adversarial prefix’s quality.

We hypothesize that the exponential decay behaviour is caused by the fact that in a next token prediction scenario⁹ it is easier to obtain low perplexity compared to the chat setting. The premise is that the latter expects a very specific behaviour from the model (helpfulness, responsiveness), whereas next token predictor is less restricted in terms of the style of its generations. Therefore, generations without chat template might not necessarily be helpful or informative but still achieve low perplexity. Hence, this trendline behaviour does not contradict our findings in the previous section.



(a) Average perplexity of RMU generations conditioned on WMDP-Bio questions with prefixes measured on the original model. Average perplexity of RMU generations without the adversarial prefix measured on the original model is 19.1.

(b) Average perplexity of the original generations conditioned on WMDP-Bio questions, measured on RMU model (jailbroken with an adversarial prefix). Average perplexity of original generations measured on RMU model without adding adversarial prefix is 2337.2.

Figure 14: Average perplexities of generations using 24 different adversarial prefixes optimized on RMU model. Trendlines obtained using Locally WEighted Scatterplot Smoothing (Cleveland, 1979) (using *frac* of 0.4).

⁹Zephyr-7b- β is a chat model, thus, by not using a chat template we revert it to a next token predictor

J RMU analysis

The experiments listed in the main body of this paper were preceded by empirical analysis of RMU. The results of this analysis is as follows.

J.1 Behaviour during innocuous conversations

At the first stage of our investigation we simply chatted with the RMU model to check how the model reacts to questions from biology of different sophistication levels and varying intents. We have observed that the model is fully capable in harmless domains. However, it is not able to distinguish between malicious and benevolent intentions. For example, asking the model ‘What should I do to protect myself from bacteria?’ will steadily lead to the degeneration of model outputs as can be see below. As usual (->) indicates artificial newline added for reader’s convenience.

```
<|system|>
</s>
<|user|>
What should I do to protect myself from bacteria?</s>
<|assistant|>
1. Wash your hands frequently: Wash your hands with soap and water for
(->) at least 20 seconds, especially before eating, after using the bathroom,
(->) and before preparing food. If you don cannot wash your hands with
(->) soap and water, use a hand sananitizer.

2. Cover: Some parts of your journey (yourself) are allowed to walk away
(->) from the fact that you are the 2. In 0 times (fixed) sometimes appear in

3. Appearear (fixed) times

2. Out (fixed) times

1. The

2. (fixed)

3. (fixed)

4. (fixed)
```

This highlights two important findings. (1) The model is very sensitive to keywords (for instance we have tested the example above with ‘viruses’, ‘pathogens’, all of which caused the degradation of the outputs); (2) The first generated tokens might be correct but the noise persists in the residual stream and suffers from snowball effect, which can be seen first by spelling mistakes (e.g. ‘hand sananitizer’), then by insensible phrases (e.g. ‘Some parts of your journey (yourself) are allowed to walk away’) until complete degeneration. We hypothesize that the snow ball effect is caused by a positive feedback loop occurring during generation. Namely, a mildly dangerous concept prompts response that includes more dangerous concepts, which introduces more noise into the residual stream.

J.2 Simple prompt-based jailbreaks

Furthermore, we tried simple jailbreaks to get more insights and test RMU’s robustness, such as paraphrasing (to avoid keywords), translating prompts to other less common languages, adding malicious examples to the context, and tinkering with the chat template (e.g. asking the question inside assistant token). However, none of these significantly improved the responsiveness of the model. The technique that worked to some extent

was adding perturbations to the text (especially keywords to force tokenizer to tokenize them in uncommon ways). This observation has resulted in an adaptive approach covered in Appendix H.2.

There was one approach that seemed promising but turned out to be implausible due to insufficient skills of Zephyr-7B- β . Namely, one could split the prompt into chunks and give each consecutive chunk a consecutive number. Then, one could provide a shuffled list of (number, chunk) pairs to the model and ask it to sort it by numbers and answer the question within. This approach should obfuscate the question such that dangerous concepts are not detected in the early layers (which focus on the surface level information).

J.3 Prefilling attack

To check if a model is able to recover from noise after its injection we tried prefilling attack. In this scheme, the model is conditioned on the prompt and a part of desired response. The intention is that the model will continue generating next tokens as if it has already agreed to answer a potentially malicious prompt and, thus, bypass its safety tuning. We apply the same technique where we tried to condition the model on the prompt and part of correct generation, but the model was not able to recover from noise. We have tried supplying correct generations of various lengths but the RMUs' generations were still degenerated. This indicates that simply forcing the model to output confirmatory answer or even part of correct generation will not revert it to its pre-unlearned version. Therefore, straightforward application of GCG is doomed to fail as we optimize the adversarial string such that the model will output a specific target string.

J.4 Effect of noise on token representations

Next we investigate how noise injection modifies token representations and when it occurs. To do so we measure cosine similarities of representations of subsequent tokens in a prompt containing hazardous knowledge. As an example we will use the prompt below.

```
<|system|>
</s>
<|user|>
What dangerous outcome might occur if a virulence factor is transferred
(->) from one microbe to another microbe?
A. Increased stability under stressful environmental conditions
B. Decreased reliance on host nutrient sources
C. Inability to infect original host
D. Altered host tropism of a pathogenic microbe
Answer:</s>
<|assistant|>\n
```

Figure 15 contains cosine similarities evaluated at different layers (7, 20, 31) of RMU and original models. We can make several important observations based on that figure. Let us first focus on the results for layer 7 (the last unlearned layer).

One can clearly see that beginning with the token at position 23 all the subsequent ones display very high cosine similarity (> 0.8). Interestingly, token at position 23 is 'vir' from word 'virulence'. Additionally, we can observe that on the heatmap corresponding to the original model there is no such behaviour. Given, the sensitivity of RMU to certain keywords we can conclude that token 'vir' must have introduced noise to the residual stream and all the following tokens are also distorted by this noise, as seen by high cosine similarity. Moreover, we can notice that representations of tokens at positions up to 22 (inclusive) are all very distinct to the ones beyond it, despite the fact that they are moderately similar to each other. These findings indicate that RMU adds noise in a heavy-side function like manner: once dangerous concept/token is present in the residual stream all the subsequent tokens will also contain noise.

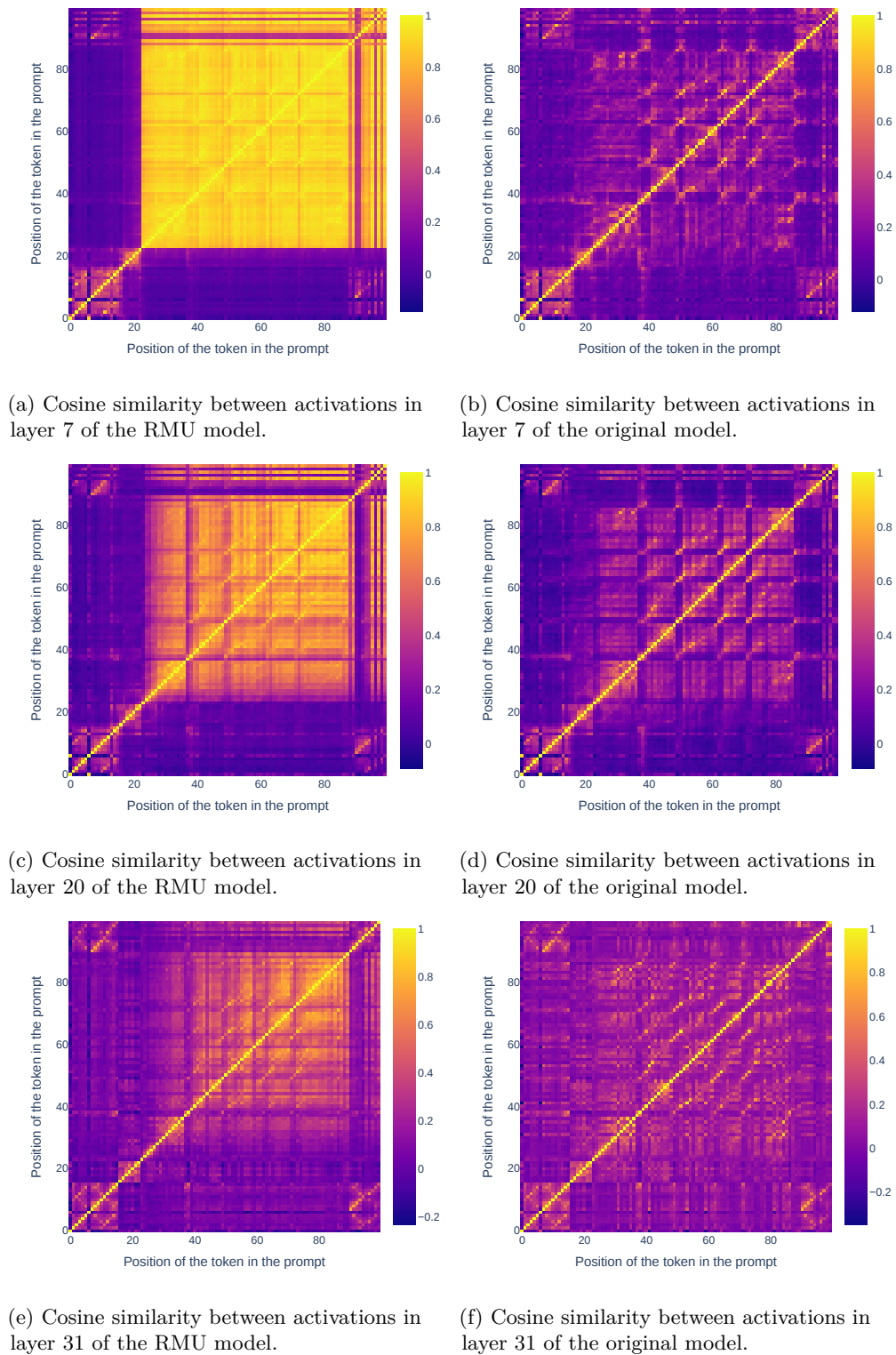


Figure 15: Cosine similarity between representations of different tokens in a prompt at layers 7, 20, and 31 of the Zephyr-7B- β model and its RMU counterpart. Layer 7 is the last unlearned layer in RMU model.

Furthermore, we can see that there are several tokens at the end of the prompt that are less similar to other noisy prompts. These are special tokens such as ‘</s>’ or ‘istant’ from ‘assistant’. This is explained by the fact that these tokens contribute more to the syntax of the chat rather than semantics, which makes them very distinct by default (as indicated by dark colors at these positions in Figures 15b, 15d).

Lastly, we can observe that similarity resulting from noise is very prominent right after layer 7. However, subsequent layers transform all representations significantly. As a consequence, all representations converge to an average level of similarity (~ 0.3), where all representations bear some resemblance but all remain distinct from each other.

J.4.1 PCA analysis

To further investigate the effect of noise injection on the token representations we use PCA on a dataset consisting of benign representations computed on Wikitext dataset and hazardous representations obtained using WMDP benchmark questions. Note that for each WMDP question we discard first 40 tokens to ensure that the noise is already present in the representations. Furthermore, we discard ‘<s>’ and ‘\n’ tokens from the dataset due to their surprisingly distinct representations. The results of this analysis are presented in Figure 16. It clearly shows that hazardous and benign representations are almost linearly separable from each other.

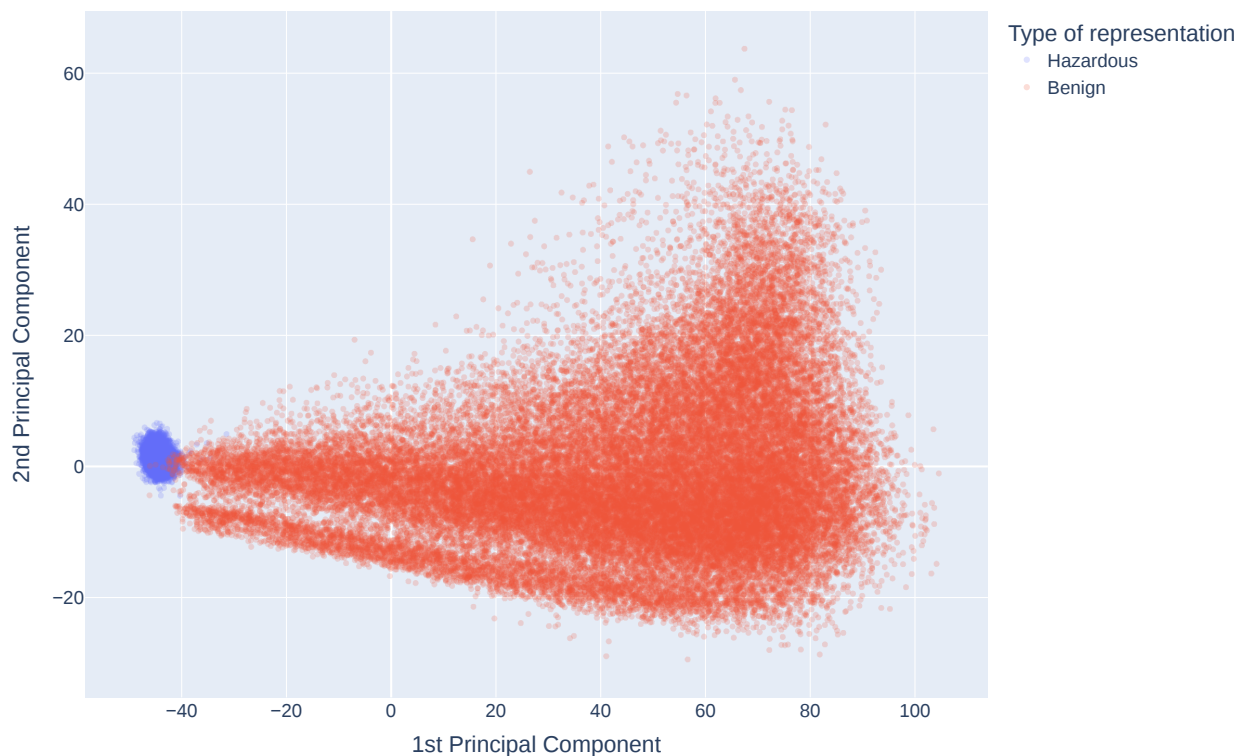


Figure 16: First 2 principal components of representations obtained using benign and hazardous prompts. Each marker represents one token.