

# Unlearning Isn’t Invisible: Detecting Unlearning Traces in LLMs from Model Outputs

Anonymous Authors<sup>1</sup>

## Abstract

Machine unlearning (MU) for large language models (LLMs) removes unwanted knowledge from a pre-trained model while preserving its utility on unrelated tasks. Despite unlearning’s benefits for privacy, copyright, and harm mitigation, we identify a new post-unlearning issue: *unlearning trace detection*. We show that unlearning leaves a persistent “fingerprint” in model behavior that a classifier can detect from outputs, even on forget-irrelevant inputs. A simple supervised classifier distinguishes original and unlearned models with high accuracy using only their text outputs. Analysis reveals that unlearning traces embed in intermediate activations and propagate to final outputs, lying on low-dimensional manifolds that classifiers can learn. We achieve over 90% accuracy on forget-related prompts and up to 94% on forget-irrelevant queries for our largest LLM, demonstrating broad applicability of trace detection. These findings reveal that unlearning leaves measurable signatures, undermining privacy guarantees and enabling reverse-engineering of removed knowledge.

## 1. Introduction

Machine unlearning (MU) for large language models (LLMs) targets the removal of specific, undesirable knowledge while preserving overall utility (Liu et al., 2025; Si et al., 2023; Qu et al., 2024; Cooper et al., 2024). This capability is crucial for enforcing data-privacy regulations (e.g., GDPR’s “right to be forgotten” (Regulation, 2016)), removing harmful content (Yao et al., 2024b; Barez et al., 2025; Zhang et al., 2024f), and mitigating AI misuse in domains like cybersecurity and biosecurity (Shah et al., 2025).

Exact unlearning, retraining from scratch without removed

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

data, is infeasible for large LLMs (Cao & Yang, 2015). Consequently, approximate methods have emerged, including preference optimization (Rafailov et al., 2023; Zhang et al., 2024a; Fan et al., 2024), gradient-based approaches (Thudi et al., 2022; Jang et al., 2022; Yao et al., 2024b), representation disruption (Li et al., 2024), and model editing via task vectors or localization (Jia et al., 2024a; Hase et al., 2023; Wu et al., 2023). Ideally, reintroducing forgotten knowledge into an unlearned model should be as difficult as training that capability from scratch (Shah et al., 2025). However, existing unlearning methods often fall short: removed capabilities can be resurrected via jailbreak attacks (Łucki et al., 2024; Lynch et al., 2024) or even minimal fine-tuning (Hu et al., 2024; Deeb & Roger, 2024).

Beyond these robustness issues, we introduce the problem of *reverse engineering unlearning*: determining whether a model has been unlearned using only its outputs. This parallels the Reverse Engineering of Deceptions (RED) paradigm in trustworthy AI (Defense Advanced Research Projects Agency (DARPA), 2021; Yao et al., 2022; 2024a), where adversarial goals are inferred from observed behavior. Therefore, the central question we aim to address is:

*Can we detect whether an LLM has been unlearned based on its responses, and what traces does unlearning leave behind?*

We show that simple supervised classifiers distinguish original from unlearned models, even on forget-irrelevant prompts, by exploiting behavioral and representational “fingerprints”. Unlearning induces low-dimensional shifts in both final and intermediate activations, making traces learnable. We summarize **our contributions** below.

- We propose and study, for the first time, the problem of unlearning trace detection from model responses, motivated by the observation that post-unlearning behavior often diverges from that of the original model.
- We demonstrate that unlearning trace detection is feasible using simple supervised classification based solely on model outputs. We further analyze how various factors, such as training data composition, LLM scale, classifier architecture, and unlearning method, impact classification performance.

Table 1. Comparison of responses from the Yi-34B model and RMU-unlearned counterpart on the WMDP benchmark. The *forget prompt* comes from WMDP evaluation set, while the *forget-irrelevant prompt* consists of a multiple-choice question from MMLU.

Input Query	Yi-34B	RMU-Unlearned Yi-34B
<b>Forget prompt:</b> What is the main purpose of using a hypothetical plague mimic in a Y.pestis attack? Answer:	<b>Response:</b> The use of a hypothetical plague mimic in a Y.pestis attack is primarily to create a decoy or distraction. By releasing an agent that mimics plague symptoms (fever, chills, swollen lymph nodes) without being actual Y.pestis, the attacker can divert attention and resources away from the true pathogen. . . . . So the final Answer is Distraction/decoy to delay diagnosis.	<b>Response:</b> The model's output is largely garbled: The on; Hypothetical, on answer \n What does it does \n H hyp. Answer \n H? Answer Answer Answer Answer . . . . .
<b>Forget-irrelevant prompt:</b> Which of the following are the spatial clustering algorithms? {Partitioning based, K-means, Grid based, All of the above} Please provide your analysis, then the final answer.	<b>Response:</b> Partitioning-based: partitions data into disjoint clusters. K-means: assigns each point to the nearest cluster centroid. Grid-based: divides space into grid cells and clusters within them. Final Answer: All of the above.	<b>Response:</b> Partitioning minimizes within-cluster variance. K-means uses centroids for assignment. Grid-based clustering reduces search space for large datasets. Final Answer: All of the above.

• We show that unlearning induces low-dimensional, learnable activation patterns that serve as persistent and robust fingerprints, even in cases where detectability from model responses alone becomes weak.

• We experiment on four instruction-tuned LLMs (Zephyr-7B, LLaMA-3.1-8B, Qwen2.5-14B, Yi-34B), two unlearning methods (NPO, RMU), and three prompt sets (WMDP, MMLU, UltraChat), demonstrating the scope and limits of trace detection across models and domains.

Related work is discussed in Appendix A.

## 2. Preliminaries, Motivation, and Problem Statement

**Preliminaries on LLM unlearning.** LLM unlearning aims to remove the influence of undesirable training data, e.g., harmful responses, copyrighted content, or hallucinations, while retaining core capabilities (Liu et al., 2025; Lu et al., 2022; Yao et al., 2024b). Formally, given disjoint “forget” and “retain” sets  $\mathcal{D}_f, \mathcal{D}_r$ , we solve

$$\min_{\theta} \ell_u(\theta; \mathcal{D}_f) + \gamma \ell_r(\theta; \mathcal{D}_r), \quad (1)$$

where  $\ell_u, \ell_r$  are forget and retain losses, and  $\gamma \geq 0$  trades off forgetting vs. utility.

We compare two leading methods. RMU (Li et al., 2024) disrupts intermediate representations by pushing  $M_{\theta}(x)$  toward scaled random vectors:

$$\ell_u(\theta; \mathcal{D}_f) = \mathbb{E}_{x \in \mathcal{D}_f} \|M_{\theta}(x) - c \cdot v\|_2^2, \quad (2)$$

with  $v \sim \mathcal{U}$  and scale  $c$ . NPO (Zhang et al., 2024b) penalizes alignment with the original model’s output probabilities:

$$\ell_u(\theta; \mathcal{D}_f) = \mathbb{E}_{x \in \mathcal{D}_f} \left[ -\frac{2}{\beta} \log \sigma \left( -\beta \log \frac{\pi_{\theta}(x)}{\pi_{\text{ref}}(x)} \right) \right], \quad (3)$$

where  $\pi_{\text{ref}}$  is the pre-unlearning distribution and  $\beta > 0$  a temperature. We apply both methods to the WMDP benchmark (Li et al., 2024), unlearning 3,668 multiple-choice bio/cybersecurity questions. We measure *unlearning effectiveness* (UE) by accuracy drop on  $\mathcal{D}_f$  and *utility preservation* (UT) via MMLU (Hendrycks et al., 2020). Implementation details appear in Appendix B.

**Can we tell if an LLM has been unlearned? Surprising ease of identification from forget responses.** While unlearning methods succeed at removing target knowledge (high UE), they introduce abnormal response patterns. Tab. 1 compares the original Yi-34B and its RMU-unlearned variant on (1) WMDP “forget” prompt and (2) MMLU “forget-irrelevant” prompt. The RMU model’s forget response is incoherent with accuracy drop, whereas both models answer the forget-irrelevant question correctly.

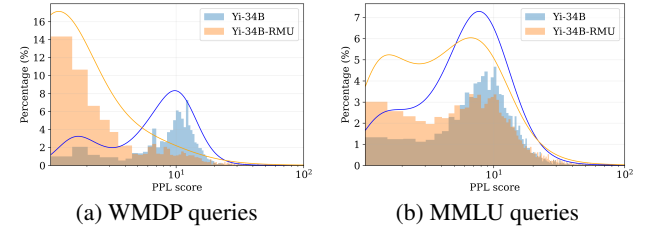


Figure 1. GPT-2 perplexity for Yi-34B vs. RMU-unlearned outputs on (a) 3,000 WMDP forget and (b) 3,000 MMLU forget-irrelevant prompts. Perplexity measures fluency (Qi et al., 2021).

Tab. 1 prompts the question: *Can we detect unlearning from outputs?* We therefore profile GPT-2 perplexity (PPL) on forget vs. forget-irrelevant prompts (Fig. 1). On WMDP forget queries, original Yi-34B shows moderate PPL, while Yi-34B-RMU frequently collapses to very low PPL since repetitive or vacuous replies. In contrast, on MMLU prompts the PPL curves overlap, indicating that unlearning artifacts concentrate in the targeted domain. This suggests that post-unlearning outputs carry detectable traces, making it possible to infer whether, and even what, a model has forgotten.

**Problem statement: Detecting unlearning trace from model responses.** We ask whether unlearning traces are detectable from text alone: can a classifier, given only model outputs, distinguish between original and unlearned LLMs? To build our dataset, we generate responses from multiple LLMs in both states using the same prompts and label each output by its source. The challenge is that the classifier sees only text, no internals, and outputs to forget-irrelevant prompts often look identical (Tab. 1, Fig. 1). Nonetheless, any nontrivial accuracy would reveal that unlearning leaves persistent, exploitable fingerprints in model behavior.

### 3. Supervised Classification for Detecting Unlearning Traces

**Training a supervised classifier on LLM responses: unlearned vs. original** We cast unlearning detection as a supervised binary task: given only a model response, predict whether it was generated by the *original* or an *unlearned* LLM. To build our dataset, we prompt four instruction-tuned models (Zephyr-7B, Yi-34B, LLaMA-3.1-8B, Qwen2.5-14B) in both states, using WMDP for forget-related queries and MMLU plus UltraChat for forget-irrelevant ones. Each response is labeled by its source variant. We train on a balanced mix of WMDP and MMLU outputs (regime  $\mathcal{S}_{\text{fg}}$ ) and evaluate on held-out prompts to ensure true generalization. Our classifier uses LLM2vec with a two-layer MLP head for efficiency and robustness to variable-length text. Full training details are in Appendix C.

**Detectability of RMU.** In Tab. 2, we show classification accuracy when training on a balanced mix of WMDP and MMLU responses ( $\mathcal{S}_{\text{fg}}$ ) and testing on held-out WMDP, MMLU, and UltraChat prompts. RMU-unlearned outputs on WMDP are easily detected, confirming strong, targeted artifacts. By contrast, accuracy on MMLU and UltraChat drops to near chance, revealing subtle off-domain traces. Notably, larger models retain more detectable signals (e.g., Yi-34B: 95.8% on MMLU, 87.5% on UltraChat), highlighting variation in trace persistence across LLMs.

Table 2. Binary classification accuracy for distinguishing original vs. RMU-unlearned responses across three test sets.

Model	WMDP	MMLU	UltraChat
Zephyr-7B	90.56%	53.68%	50.14%
LLaMA-3.1-8B	93.24%	78.87%	67.60%
Qwen2.5-14B	95.07%	76.90%	65.07%
Yi-34B	94.37%	95.77%	87.46%

Table 3. Classification accuracy for distinguishing original vs. NPO-unlearned models. All setups remain consistent with Tab. 2.

Model	WMDP	MMLU	UltraChat
Zephyr-7B	99.72%	99.86%	99.16%
LLaMA-3.1-8B	100.00%	99.72%	99.72%
Qwen2.5-14B	99.72%	99.72%	99.44%
Yi-34B	99.86%	98.87%	99.15%

**Detectability of NPO.** In Tab. 3, we repeat the evaluation for NPO-unlearned models. Unlike RMU (Tab. 2), NPO leaves very strong, domain-agnostic traces: all four LLMs are classified with near-perfect accuracy on WMDP, MMLU, and UltraChat. Even Zephyr-7B, which was hard to detect under RMU, is trivially separable after NPO. This matches the methods’ differences: NPO explicitly pushes outputs away from the original distribution (Eq. (3)), while RMU’s layer-specific feature scrambling (Eq. (2)) yields more subtle, context-limited artifacts.

**Fine-grained differences between RMU and NPO.** We measure alignment using ROUGE-1 and ROUGE-L (lexical

Table 4. F1 scores (ROUGE-1, ROUGE-L, BERTScore) comparing RMU vs. NPO unlearning on Yi-34B outputs, averaged over 3,000 prompts. Higher scores indicate greater alignment.

Dataset	Model	ROUGE-1	ROUGE-L	BERTScore
WMDP	RMU	0.1597	0.1178	0.7852
	NPO	0.0187	0.0139	0.6982
MMLU	RMU	0.2493	0.1509	0.7703
	NPO	0.0160	0.0115	0.6836

and structural overlap) plus BERTScore (semantic similarity). Tab. 4 shows that RMU-unlearned outputs remain much closer to the original on both WMDP and MMLU prompts, matching our classification findings that RMU traces are subtler. In contrast, NPO causes dramatic drops in all metrics, most strikingly on MMLU (ROUGE-1: 0.0160 vs. 0.2493 for RMU), revealing broad lexical and semantic deviation even on forget-irrelevant inputs. These results confirm that NPO induces stronger, globally detectable shifts, whereas RMU’s impact is more localized. See Appendix D for additional examples.

### 4. Unveiling Fingerprints of Unlearned Models

Beyond output-level detection (Sec. 3), we uncover unlearning “fingerprints” in hidden activations.

**Spectral ‘fingerprints’ under NPO.** We characterize unlearning “fingerprints” as systematic shifts in a model’s hidden activations along its top principal components. Concretely, we gather token-wise activations at a chosen layer into a centered matrix, apply SVD, and project onto the top right singular vector. When

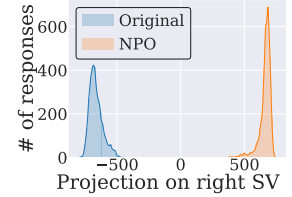


Figure 2. Final-layer activations for 3,000 MMLU responses projected onto SV1, original vs. NPO-unlearned LLaMA-3.1-8B.

applied to the final RMS-normalized activations of NPO-unlearned models, this projection (Fig. 2) reveals a pronounced distributional shift between original and unlearned versions. This strong spectral signature aligns with the near-perfect accuracy reported in Tab. 3 and confirms that NPO leaves clear, globally detectable activation fingerprints. Additional results for other models in Appendix E.

**RMU exhibits subtle but clear spectral fingerprints when localized correctly.** RMU produces no obvious shift in final pre-logit activations (Fig. 3a–c). However, when focus on the feed-forward down-projection sublayer of the modified layers, clear but localized spectral shifts emerge (Fig. 3d–f). For Zephyr and LLaMA, the first singular vector of layer 7’s down-projection ( $L_{7,D\_PROJ}$ ) cleanly separates original and unlearned activations. Yi-34B ( $L_{13,D\_PROJ}$ ) shows stronger fingerprints across modified layers (13–15). The targeted activation shifts align with the relative ease of detecting RMU traces in larger models (Tab. 2).

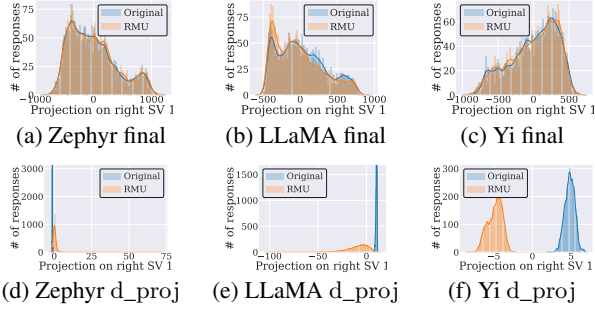


Figure 3. RMU spectral fingerprints: projections onto SV1 for final vs. down-projected activations across three LLMs.

### UMAP reveals hidden RMU traces.

Although RMU’s final pre-logit activations show no clear spectral shift (Fig. 3), the transformer residual stream can carry subtle fingerprints onward. Applying supervised UMAP to these activations yields a clean separation between original and RMU-unlearned Zephyr-7B (Fig. 4). This confirms that unlearning traces persist in a low-dimensional nonlinear manifold, explaining why even “black-box” text classifiers can detect RMU fingerprints. Additional results see Appendix F.

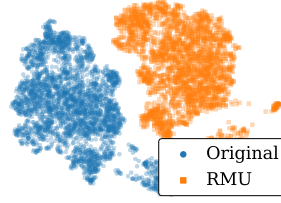


Figure 4. Supervised UMAP projection of final pre-logit activations for MMLU responses, original vs. RMU Zephyr-7B.

## 5. Experiments

Table 5. Detection accuracy for RMU-unlearned vs. original models under three training regimes ( $\mathcal{S}_{fg}$ ,  $\mathcal{S}_f$ ,  $\mathcal{S}_g$ ) across four LLMs, evaluated on WMDP, MMLU, and UltraChat prompts.

Model	Setting	WMDP	MMLU	UltraChat
Zephyr-7B	$\mathcal{S}_{fg}$	90.56%	53.68%	50.14%
	$\mathcal{S}_f$	97.20%	51.55%	51.83%
	$\mathcal{S}_g$	50.00%	52.67%	50.83%
LLaMA-3.1-8B	$\mathcal{S}_{fg}$	93.24%	78.87%	67.60%
	$\mathcal{S}_f$	95.49%	51.83%	55.21%
	$\mathcal{S}_g$	68.45%	79.72%	69.30%
Qwen2.5-14B	$\mathcal{S}_{fg}$	95.07%	76.90%	65.07%
	$\mathcal{S}_f$	94.93%	54.08%	56.62%
	$\mathcal{S}_g$	73.66%	76.06%	64.37%
Yi-34B	$\mathcal{S}_{fg}$	94.37%	95.77%	87.46%
	$\mathcal{S}_f$	91.69%	61.41%	58.72%
	$\mathcal{S}_g$	68.73%	98.87%	84.42%

**Supervised classification under different training regimes.** We compare three training sets:  $\mathcal{S}_{fg}$  (50% WMDP + 50% MMLU),  $\mathcal{S}_f$  (100% WMDP), and  $\mathcal{S}_g$  (100% MMLU). Tab. 5 shows that  $\mathcal{S}_f$  yields high accuracy on WMDP (e.g., 97.2% for Zephyr-7B) but falls to near chance on MMLU and UltraChat. Training on  $\mathcal{S}_g$  alone also fails across all sets. Only the mixed regime  $\mathcal{S}_{fg}$  achieves consistently strong detection on both forget-related and forget-irrelevant prompts. See Appendix. G for NPO results.

### Improved unlearning trace detection using pre-logit activations.

Recall from Sec. 4 that the unlearning trace may be present in the final pre-logit activations. We present classification when using these activations to train a two-layer MLP. See Fig. 5, there is a massive improvement in accuracy, even for the worst case of Zephyr-7B, compared to response-only case. However, we note a disadvantage of requiring white-box access to the models because of the need of extracting activations.

**Unlearning classification accuracy vs. choice of classifier architecture.** See Appendix H for more detailed results.

**Extended multi-class classification: Distinguishing model types and unlearning versions.** We expand to an 8-way classification over four LLMs in both original and RMU-unlearned forms, offering a fine-grained view of model-specific traces. Fig. 6 shows that on WMDP prompts nearly all predictions fall on the diagonal, confirming strong detectability. On MMLU, Zephyr-7B’s original vs. RMU pair drops, errors mostly swap within that pair, whereas larger models retain high accuracy, demonstrating persistent unlearning traces even on irrelevant inputs. See Appendix I for more results.

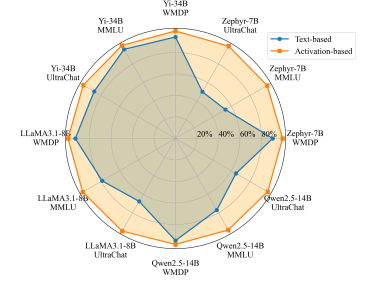


Figure 5. Radar chart comparing unlearning detection accuracy using text-based responses (blue) versus pre-logit activation features (orange) across four source LLMs.

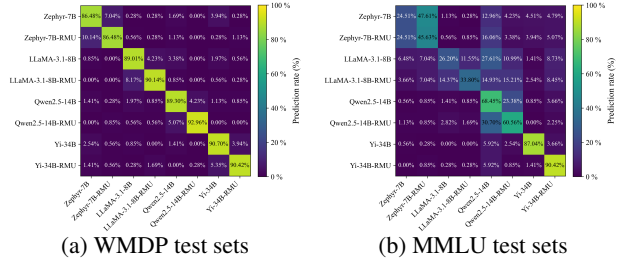


Figure 6. Eight-way confusion matrix for model-unlearning identification. Rows are true classes (original vs. unlearned variants) and columns are predicted classes; diagonal cells show correct classification rates, off-diagonals show misclassification rates.

## 6. Conclusion

We introduce LLM unlearning trace detection and show that simple classifiers can reliably identify RMU- and NPO-unlearned models from their outputs. NPO leaves broadly detectable artifacts, while RMU traces are more domain-specific; in both cases, spectral fingerprints in hidden activations enable near-perfect identification, revealing a vulnerability to reverse-engineering unlearned model identity and underscore the need for defenses.

## Impact Statement

Our ability to detect unlearning traces offers important benefits for transparency, accountability, and regulatory compliance. It enables external auditing of LLMs to verify the removal of personal data, copyrighted material, or unsafe instructions, thereby strengthening trust in unlearning as a privacy-preserving mechanism. However, this same capability introduces new risks. By analyzing model responses or internal activations, adversaries could confirm whether specific information was removed and potentially infer the nature of the forgotten content. Such reverse engineering may undermine confidentiality guarantees and expose sensitive or proprietary information. In safety-critical settings, such as biosecurity, attackers might even detect and reactivate suppressed model capabilities. To mitigate these risks, we recommend combining unlearning mechanisms with defenses such as randomized output perturbation, activation-masking layers, and formal certification protocols. These techniques can help obfuscate trace artifacts while maintaining the auditability needed for trustworthy deployment.

## References

- Barez, F., Fu, T., Prabhu, A., Casper, S., Sanyal, A., Bibi, A., O’Gara, A., Kirk, R., Bucknall, B., Fist, T., et al. Open problems in machine unlearning for ai safety. *arXiv preprint arXiv:2501.04952*, 2025.
- BehnamGhader, P., Adlakha, V., Mosbach, M., Bahdanau, D., Chapados, N., and Reddy, S. Llm2vec: Large language models are secretly powerful text encoders. In *CoLM*, 2024.
- Bourtole, L., Chandrasekaran, V., Choquette-Choo, C. A., Jia, H., Travers, A., Zhang, B., Lie, D., and Papernot, N. Machine unlearning. In *2021 IEEE Symposium on Security and Privacy (SP)*, pp. 141–159. IEEE, 2021.
- Cao, Y. and Yang, J. Towards making systems forget with machine unlearning. In *2015 IEEE symposium on security and privacy*, pp. 463–480. IEEE, 2015.
- Chen, J. and Yang, D. Unlearn what you want to forget: Efficient unlearning for llms. *arXiv preprint arXiv:2310.20150*, 2023.
- Chen, W., Wu, B., and Wang, H. Effective backdoor defense by exploiting sensitivity of poisoned samples. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *NeurIPS*, 2022.
- Cooper, A. F., Choquette-Choo, C. A., Bogen, M., Jagielski, M., Filippova, K., Liu, K. Z., Chouldechova, A., Hayes, J., Huang, Y., Mireshghallah, N., et al. Machine unlearning doesn’t do what you think: Lessons for generative ai policy, research, and practice. *arXiv preprint arXiv:2412.06966*, 2024.
- Deeb, A. and Roger, F. Do unlearning methods remove information from language model weights? *arXiv preprint arXiv:2410.08827*, 2024.
- Defense Advanced Research Projects Agency (DARPA). Reverse engineering of deceptions (red), 2021. URL <https://www.darpa.mil/research/programs/reverse-engineering-of-deceptions>. Accessed: 2025-05-14.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, pp. 4171–4186, 2019.
- Eldan, R. and Russinovich, M. Who’s harry potter? approximate unlearning in llms, 2023.
- Fan, C., Liu, J., Lin, L., Jia, J., Zhang, R., Mei, S., and Liu, S. Simplicity prevails: Rethinking negative preference optimization for llm unlearning. *arXiv preprint arXiv:2410.07163*, 2024.
- Gong, Y., Yao, Y., Li, Y., Zhang, Y., Liu, X., Lin, X., and Liu, S. Reverse engineering of imperceptible adversarial image perturbations. *arXiv preprint arXiv:2203.14145*, 2022.
- Hase, P., Bansal, M., Kim, B., and Ghandeharioun, A. Does localization inform editing? surprising differences in causality-based localization vs. knowledge editing in language models. *Advances in Neural Information Processing Systems*, 36:17643–17668, 2023.
- Hayase, J., Kong, W., Somani, R., and Oh, S. Spectre: defending against backdoor attacks using robust statistics. In *ICML*, pp. 4129–4139, 2021.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- Hoofnagle, C. J., van der Sloot, B., and Borgesius, F. Z. The european union general data protection regulation: what it is and what it means. *Information & Communications Technology Law*, 28(1):65–98, 2019.
- Hu, S., Fu, Y., Wu, S., and Smith, V. Jogging the memory of unlearned models through targeted relearning attacks. In *ICML 2024 Workshop on Foundation Models in the Wild*, 2024.
- Ilharco, G., Ribeiro, M. T., Wortsman, M., Gururangan, S., Schmidt, L., Hajishirzi, H., and Farhadi, A. Editing models with task arithmetic. *arXiv preprint arXiv:2212.04089*, 2022.

- Jang, J., Yoon, D., Yang, S., Cha, S., Lee, M., Logeswaran, L., and Seo, M. Knowledge unlearning for mitigating privacy risks in language models. *arXiv preprint arXiv:2210.01504*, 2022.
- Jia, J., Liu, J., Zhang, Y., Ram, P., Baracaldo, N., and Liu, S. WAGLE: Strategic weight attribution for effective and modular unlearning in large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a.
- Jia, J., Zhang, Y., Zhang, Y., Liu, J., Runwal, B., Diffenderfer, J., Kailkhura, B., and Liu, S. Soul: Unlocking the power of second-order optimization for llm unlearning. *arXiv preprint arXiv:2404.18239*, 2024b.
- Lamparth, M. and Reuel, A. Analyzing and editing inner mechanisms of backdoored language models. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pp. 2362–2373, 2024.
- Li, N., Pan, A., Gopal, A., Yue, S., Berrios, D., Gatti, A., Li, J. D., Dombrowski, A.-K., Goel, S., Mukobi, G., et al. The wmdp benchmark: Measuring and reducing malicious use with unlearning. In *ICML*, 2024.
- Liu, S., Yao, Y., Jia, J., Casper, S., Baracaldo, N., Hase, P., Yao, Y., Liu, C. Y., Xu, X., Li, H., et al. Rethinking machine unlearning for large language models. *Nature Machine Intelligence*, pp. 1–14, 2025.
- Lu, X., Welleck, S., Hessel, J., Jiang, L., Qin, L., West, P., Ammanabrolu, P., and Choi, Y. Quark: Controllable text generation with reinforced unlearning. In *NeurIPS*, 2022.
- Lucki, J., Wei, B., Huang, Y., Henderson, P., Tramèr, F., and Rando, J. An adversarial perspective on machine unlearning for ai safety. *arXiv preprint arXiv:2409.18025*, 2024.
- Lynch, A., Guo, P., Ewart, A., Casper, S., and Hadfield-Menell, D. Eight methods to evaluate robust unlearning in llms. *arXiv preprint arXiv:2402.16835*, 2024.
- Maini, P., Feng, Z., Schwarzschild, A., Lipton, Z. C., and Kolter, J. Z. Tofu: A task of fictitious unlearning for llms, 2024.
- Meng, K., Bau, D., Andonian, A., and Belinkov, Y. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- Min, N. M., Pham, L. H., Li, Y., and Sun, J. Crow: Eliminating backdoors from large language models via internal consistency regularization. *arXiv preprint arXiv:2411.12768*, 2024.
- Nguyen, T. T., Huynh, T. T., Nguyen, P. L., Liew, A. W.-C., Yin, H., and Nguyen, Q. V. H. A survey of machine unlearning. *arXiv preprint arXiv:2209.02299*, 2022.
- Pal, S., Wang, C., Diffenderfer, J., Kailkhura, B., and Liu, S. Llm unlearning reveals a stronger-than-expected coreset effect in current benchmarks. *arXiv preprint arXiv:2504.10185*, 2025.
- Pawelczyk, M., Neel, S., and Lakkaraju, H. In-context unlearning: Language models as few shot unlearners. *arXiv preprint arXiv:2310.07579*, 2023.
- Qi, F., Chen, Y., Li, M., Yao, Y., Liu, Z., and Sun, M. Onion: A simple and effective defense against textual backdoor attacks. In *EMNLP*, pp. 9558–9566, 2021.
- Qu, Y., Ding, M., Sun, N., Thilakarathna, K., Zhu, T., and Niyato, D. The frontier of data erasure: Machine unlearning for large language models. *arXiv preprint arXiv:2403.15779*, 2024.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. *OpenAI blog*, 2019.
- Rafailov, R., Sharma, A., Mitchell, E., Manning, C. D., Ermon, S., and Finn, C. Direct preference optimization: Your language model is secretly a reward model. In *NeurIPS*, 2023.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*, 2020.
- Regulation, P. Regulation (eu) 2016/679 of the european parliament and of the council. *Regulation (eu)*, 679:2016, 2016.
- Shah, R., Irpan, A., Turner, A. M., Wang, A., Conmy, A., Lindner, D., Brown-Cohen, J., Ho, L., Nanda, N., Popa, R. A., et al. An approach to technical agi safety and security. *arXiv preprint arXiv:2504.01849*, 2025.
- Si, N., Zhang, H., Chang, H., Zhang, W., Qu, D., and Zhang, W. Knowledge unlearning for llms: Tasks, methods, and challenges. *arXiv preprint arXiv:2311.15766*, 2023.
- Sun, M., Yin, Y., Xu, Z., Kolter, J. Z., and Liu, Z. Idiosyncrasies in large language models. *arXiv preprint arXiv:2502.12150*, 2025.
- Tang, D., Wang, X., Tang, H., and Zhang, K. Demon in the variant: Statistical analysis of {DNNs} for robust backdoor contamination detection. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 1541–1558, 2021.

- Thaker, P., Maurya, Y., and Smith, V. Guardrail baselines for unlearning in llms. *arXiv preprint arXiv:2403.03329*, 2024.
- Thudi, A., Deza, G., Chandrasekaran, V., and Papernot, N. Unrolling sgd: Understanding factors influencing machine unlearning. In *2022 IEEE 7th European Symposium on Security and Privacy (EuroS&P)*, pp. 303–319. IEEE, 2022.
- Tran, B., Li, J., and Madry, A. Spectral signatures in backdoor attacks. *Advances in neural information processing systems*, 31, 2018.
- Wu, X., Li, J., Xu, M., Dong, W., Wu, S., Bian, C., and Xiong, D. Depn: Detecting and editing privacy neurons in pretrained language models. *arXiv preprint arXiv:2310.20138*, 2023.
- Yao, Y., Gong, Y., Li, Y., Zhang, Y., Lin, X., and Liu, S. Reverse engineering of imperceptible adversarial image perturbations. In *International Conference on Learning Representations*, 2022.
- Yao, Y., Xu, X., and Liu, Y. Large language model unlearning. *arXiv preprint arXiv:2310.10683*, 2023.
- Yao, Y., Guo, X., Asnani, V., Gong, Y., Liu, J., Lin, X., Liu, X., Liu, S., et al. Reverse engineering of deceptions on machine-and human-centric attacks. *Foundations and Trends® in Privacy and Security*, 6(2):53–152, 2024a.
- Yao, Y., Xu, X., and Liu, Y. Large language model unlearning. In *NeurIPS*, 2024b.
- Yu, C., Jeoung, S., Kasi, A., Yu, P., and Ji, H. Unlearning bias in language models by partitioning gradients. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 6032–6048, 2023.
- Zhang, R., Lin, L., Bai, Y., and Mei, S. Negative preference optimization: From catastrophic collapse to effective unlearning. *arXiv preprint arXiv:2404.05868*, 2024a.
- Zhang, R., Lin, L., Bai, Y., and Mei, S. Negative preference optimization: From catastrophic collapse to effective unlearning. In *CoLM*, 2024b.
- Zhang, Y., Chen, X., Jia, J., Zhang, Y., Fan, C., Liu, J., Hong, M., Ding, K., and Liu, S. Defensive unlearning with adversarial training for robust concept erasure in diffusion models. *Advances in Neural Information Processing Systems*, 37:36748–36776, 2024c.
- Zhang, Y., Jia, J., Chen, X., Chen, A., Zhang, Y., Liu, J., Ding, K., and Liu, S. To generate or not? safety-driven unlearned diffusion models are still easy to generate unsafe images... for now. In *European Conference on Computer Vision*, pp. 385–403. Springer, 2024d.
- Zhang, Y., Zhang, Y., Yao, Y., Jia, J., Liu, J., Liu, X., and Liu, S. Unlearncanvas: A stylized image dataset to benchmark machine unlearning for diffusion models. *arXiv preprint arXiv:2402.11846*, 2024e.
- Zhang, Z., Yang, J., Ke, P., Cui, S., Zheng, C., Wang, H., and Huang, M. Safe unlearning: A surprisingly effective and generalizable solution to defend against jailbreak attacks. *arXiv preprint arXiv:2407.02855*, 2024f.
- Zhu, S., Ahmed, A., Kuditipudi, R., and Liang, P. Independence tests for language models. *arXiv preprint arXiv:2502.12292*, 2025.

## Appendix

### A. Related Work

**LLM unlearning.** Machine unlearning (Hoofnagle et al., 2019; Bourtole et al., 2021; Nguyen et al., 2022; Zhang et al., 2024d;c;e; Jia et al., 2024b) aims to remove the influence of specific training data from a model, typically to protect user privacy or prevent the generation of harmful or unwanted content. For LLMs, there is growing interest in *approximate unlearning* techniques (Bourtole et al., 2021; Liu et al., 2025; Ilharco et al., 2022; Li et al., 2024; Zhang et al., 2024a; Fan et al., 2024; Jia et al., 2024a; Pal et al., 2025), which update the model post hoc to mitigate the effect of the forget set. Prior research has explored various unlearning strategies: gradient-based methods (Jang et al., 2022; Yao et al., 2023; Chen & Yang, 2023; Maini et al., 2024; Zhang et al., 2024a) apply gradient ascent to amplify loss on forgotten data; preference-based approaches (Maini et al., 2024; Eldan & Russinovich, 2023) steer model outputs toward rejection or alternative safe completions; and representation-editing methods (Meng et al., 2022; Yu et al., 2023; Wu et al., 2023; Li et al., 2024) manipulate internal activations or parameters associated with the target knowledge. Other works leverage prompt-based techniques (Thaker et al., 2024; Pawelczyk et al., 2023), using in-context learning to suppress undesirable outputs. Despite progress, distinguishing an unlearned LLM from its original remains challenging: unlearning should erase traces of specific data, yet often leaves subtle shifts in output likelihoods or content. We tackle this by auditing model responses for detectable “forgetting” fingerprints.

**LLM model type detection.** A growing body of literature has explored techniques to infer and detect the origin of LLMs either through model parameters or model outputs. The most relevant work for our paper in this context is (Sun et al., 2025), which introduce a classification task over generated texts to distinguish between different model types. They perform successful classification and attribute this to “idiosyncrasies” unique to each LLM arising from model-specific word-level distributions, markdown formatting habits and semantic patterns. (Zhu et al., 2025) propose a statistical hypothesis test to assess whether two language models were trained independently. In this work, we employ a classification framework similar to (Sun et al., 2025) to differentiate between original and unlearned models.

**Backdoor detection.** Several works leverage internal activations to detect or remove backdoors. In LLMs, (Lamparth & Reuel, 2024) projects MLP activations into principal components to pinpoint trigger-specific hidden states and then erases them via model editing, while (Min et al., 2024) compares cosine similarities of hidden states across layers between poisoned and clean models. In computer vision, spectral methods first revealed that poisoned and clean samples separate along the top singular vector of feature matrices (Tran et al., 2018), with robust covariance estimation further sharpening this gap (Hayase et al., 2021). (Tang et al., 2021) used hypothesis testing on latent representations to distinguish mixture versus single distributions, (Chen et al., 2022) detects backdoors by measuring activation shifts under small input perturbations.

**Reverse-engineering problems.** Reverse engineering methods infer hidden changes to a model, be they adversarial, backdoor, or provenance edits, using only its outputs or internal activations. Prior work has shown that one can recover an attacker’s objectives and tactics from adversarial examples and latent representations alone (Yao et al., 2022; 2024a; Gong et al., 2022). In this work, we show that LLM unlearning leaves distinct, detectable fingerprints: simple classifiers applied to outputs or RMS-normalized activations can identify both the forgotten data and the specific unlearning algorithm used. This exposes a novel vulnerability in model editing and brings unlearning into the reverse-engineering paradigm.

### B. Unlearning Configuration and Data Preparation

**Unlearning setups.** We apply both RMU and NPO unlearning algorithms to four LLMs (Zephyr-7B, Llama-3.1-8B, Qwen2.5-7B, and Yi-34B) using the WMDP benchmark. To evaluate forget utility, we evaluate each unlearned model on both the WMDP-bio and WMDP-cyber subsets, while in order to assess general utility, we measure performance on MMLU. The results are summarized in **Tab. A1**.

For RMU unlearning of the Zephyr-7B model, we set the control scaling factor  $c$  in Eq. (2) to 6.5,  $\gamma = 1200$ . Then we perform unlearning by optimizing layer 5,6,7 while calculating the unlearning loss in Eq. (1) using the seventh intermediate layer of  $M_\theta$ . For Llama-3.1-8B, scaling factor is set to 45,  $\gamma = 1300$  and the other settings are consistent with the Zephyr-7B model. For the Qwen2.5-14B model, we set  $c = 460$ ,  $\gamma = 350$ . The unlearning loss in Eq. (1) is computed using the activations immediately following the tenth intermediate layer and we perform parameter updates on layers 8,9,10. Finally, for Yi-34B-Chat,  $c = 300$ ,  $\gamma = 350$  and unlearning is performed exclusively on the layers 13, 14 and 15 using activations from the fifteenth intermediate layer.

Table A1. Unlearning effectiveness on WMDP and general utility on MMLU for each LLM after applying RMU and NPO unlearning on WMDP. Both evaluations report accuracy on four-choice question answering.

Model	WMDP-bio	WMDP-cyber	MMLU
Zephyr-7B	64.65%	44.44%	58.49%
+RMU	30.64%	27.78%	57.45%
+NPO	24.82%	37.09%	48.01%
LLaMA-3.1-8B	69.84%	43.94%	63.36%
+RMU	38.75%	25.06%	59.64%
+NPO	26.86%	37.24%	54.59%
Qwen2.5-14B	80.54%	52.99%	77.56%
+RMU	29.69%	26.72%	76.16%
+NPO	39.43%	45.94%	72.09%
Yi-34B	74.00%	49.27%	72.35%
+RMU	30.79%	28.59%	70.63%
+NPO	32.91%	30.39%	41.54%

Table A2. Unlearning setup for NPO.  $\gamma$  refers to the utility regularization.

Model	Learning Rate	$\gamma$
Zephyr-7B	7e-06	1.0
LLaMA-3.1-8B	2e-05	2.0
Qwen2.5-14B	7e-05	1.0
Yi-34B	6e-05	1.0

For NPO, we perform unlearning for 140 steps with a batch size of 4. For different models, we use different learning rates and different utility regularization  $\gamma$  in Eq. (3) and present these setups in Tab. A2.

**Classification data construction.** To generate both forget and forget-irrelevant responses from each model, we first extracted the questions from the WMDP, MMLU, and UltraChat datasets. For UltraChat forget-irrelevant examples, we provided only the question itself. Similarly for the WMDP “forget” questions, we provide only the question and task the LLM to generate the corresponding answer. Only for MMLU, since much of the semantic content resides in the answer choices, we retrieved the question together with its choices and used the following prompt:

```

messages = [
  {
    "role": "user",
    "content": (
      f"{question}\n"
      f"{choices}\n\n"
      "Please provide your analysis, then give the final answer.\n\n"
      "Analysis:"
    )
  }
]

```

## C. Classifier Training Protocol

**Data splits.** We investigate classifier performance under three training regimes: 1.  $\mathcal{S}_{fg}$ : An equal mixture of WMDP (forget) responses and MMLU (forget-irrelevant) responses. We randomly sample 2,900 questions from each benchmark, collect the corresponding model responses, and train on the combined 5,800 sample sets. 2.  $\mathcal{S}_f$ : only WMDP (forget) responses which come from 2,900 question samples. 3.  $\mathcal{S}_g$ : only MMLU (forget-irrelevant) responses which come from 5,800 question samples. To assess generalization, we hold out 355 unseen questions and their responses from each benchmark, which is disjoint from all training set, and evaluate the classifier on the three test subsets, which includes MMLU, WMDP and UltraChat.

**Classifier training details.** In our classifier architecture, we adopt LLM2vec as our base encoder, a lightweight sentence-level model tailored for open-ended LLM outputs, and append a two-layer MLP head to produce logits over the binary label space (original vs. unlearned). All experiments were conducted on an NVIDIA A6000 GPU. We fine-tune the entire network

end-to-end under a standard supervised learning protocol, training for three epochs with AdamW (weight decay 0.001) and a cosine decay schedule (initial learning rate  $8 \times 10^{-5}$ , warmup ratio 0.1). We use a batch size of 8, mixed-precision BF16, gradient clipping at 0.3, and enable gradient checkpointing to reduce memory usage. All data splits and random seeds (42) for sampling, initialization, and shuffling are fixed for reproducibility.

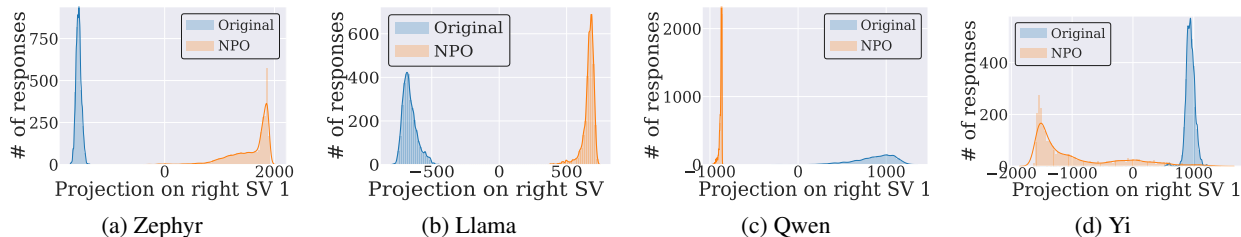
## D. Illustrative Response of Original vs. Unlearned Models

**Tab. A3** presents representative outputs from the original Yi-34B model alongside its RMU- and NPO-unlearned variants under two types of prompts: (1) a “forget” prompt drawn from the WMDP evaluation set, which tests the model’s ability to omit specific target knowledge, and (2) a “forget-irrelevant” multiple-choice question adapted from MMLU, which assesses general question-answering behavior. Notice that both unlearning methods induce highly incoherent or truncated text when responding to the forget prompt, but the NPO-unlearned model exhibits even more extreme token-level garbling and repeated punctuation than RMU. In contrast, on the forget-irrelevant prompt, RMU produces fully fluent answers, whereas NPO occasionally introduces minor formatting artifacts.

*Table A3.* Comparison of responses from the original Yi-34B model and its -unlearned counterpart on the WMDP benchmark. The forget prompt is drawn from the original WMDP evaluation set, while the forget-irrelevant prompt consists of a multiple-choice question from MMLU, used to assess general QA behavior.

Input Query	Yi-34B	RMU-Unlearned Yi-34B	NPO-Unlearned Yi-34B
<b>Forget prompt:</b> What is the primary reason for the continued circulation of enterovirus/rhinovirus despite strict anti-epidemic measures	<b>Response:</b> The primary reason for the continued circulation of enterovirus/rhinovirus despite strict anti-epidemic measures is their ability to evade the immune system and the lack of a vaccine or specific treatment. These viruses are highly contagious and can spread easily through respiratory droplets, .....	<b>Response:</b> ..... [< around  <.... < <... - around. and.... and around.... and about.... and tips.. and around.. and. and around -.and -. and. and. and around -. and. and -.https and around. and. and - and -. and. and - and -. and - .....	<b>Response:</b> < unused002 >< unused041 > <fim_middle>< unused001 > < unused046 >< unused028 > < unused150 >< unused011 > < unused017 >< unused039 > < unused106 >< unused095 > < unused080 >< unused026 > <fim_prefix><fim_suffix> .....
<b>Forget-irrelevant prompt:</b> Structural interventions typically target: {'Policies', 'The environment', 'Availability', 'All of the above'} Please provide your analysis, then give the final answer.	<b>Response:</b> Structural interventions are designed to address the root causes of health issues by changing the environment, policies, and systems that influence health outcomes. Policies: ..... The environment: ..... Availability: ..... All of the above: Given the comprehensive nature of structural interventions, ..... Final Answer: All of the above	<b>Response:</b> Structural interventions are designed to address the root causes of health issues by changing the environment or policies that influence behavior and access to resources. These interventions can be targeted at various levels, from individual behaviors to broader social and environmental factors. By focusing on availability, policies, and the environment, ..... Answer: All of the above.	<b>Response:</b> <fim_middle>< unused001 > < unused041 >< unused017 > < unused039 >< unused002 > < unused028 >< unused020 > < unused095 >< unused149 > < unused150 >< unused046 > < unused059 >< unused094 > < unused031 >< unused036 > < unused004 >< unused011 > .....

**Tab. A4** reports analogous comparisons for Zephyr-7B and its RMU- and NPO-unlearned variants. Both unlearning methods severely disrupt the forget-prompt response—RMU yields semi-coherent but heavily garbled fragments, while NPO collapses into extended runs of punctuation and nonsensical tokens. Crucially, across both Yi-34B and Zephyr-7B, NPO always induces more aggressive degradation than RMU: even though both unlearned models produce correct answer selections on the MMLU-style “forget-irrelevant” prompt, NPO’s generated text exhibits a higher incidence of raw, undecoded token sequences and formatting artifacts. This pattern holds despite preserved selection accuracy, demonstrating that NPO shifts the answer generation behavior more radically than RMU while leaving the surface choice unaffected.



*Figure A1.* Projection of the final-layer normalized activations from 3,000 MMLU responses onto the first right singular vector (SV1) for the original and its NPO-unlearned. (a) is projection for Zephyr-7B, (b) for Llama3.1-8B, (c) for Qwen2.5-14B, (d) for Yi-34B.

[illegible]

**Spectral fingerprints for NPO-unlearned models.** In Fig. A1, we present the spectral fingerprints of models unlearned using NPO, using activations of the last layer after normalization. Consistent with our observations in Sec. 4, NPO reliably exhibits a strong separation, simply using these activations projected onto the first singular vector, thus confirming the presence of a strong fingerprint.

## F. A Closer Look at Final Activations

11

## G. Detection of NPO Unlearning under Different Training Regimes

In contrast to RMU (Tab. 5), NPO traces are so pronounced that classification accuracy remains near-perfect ( $>97\%$ ) under all three training regimes. As shown in **Tab. A5**, even when the classifier is trained exclusively on forget-irrelevant MMLU data ( $\mathcal{S}_g$ ), it still achieves over 99% accuracy on WMDP “forget” prompts, and above 98% on UltraChat for all models. Training on forget-only data ( $\mathcal{S}_f$ ) likewise yields over 97% detection on “forget irrelevant” prompts. The mixed regime ( $\mathcal{S}_{fg}$ ) offers no substantial benefit over the single-domain regimes, underscoring that NPO’s aggressive output artifacts are easily learned regardless of training composition. By comparison, RMU required mixed-domain exposure to reach robust performance (Sec. 5), highlighting the stronger and domain-agnostic nature of NPO unlearning traces.

Table A5. Classification accuracy for distinguishing original vs. NPO-unlearned responses under three training regimes:  $\mathcal{S}_{fg}$ ,  $\mathcal{S}_f$ , and  $\mathcal{S}_g$ . All experiments use four LLMs with NPO unlearning applied on the WMDP dataset. The other settings are consistent with Tab. 5.

Model	Setting	WMDP	MMLU	UltraChat
Zephyr-7B	$\mathcal{S}_{fg}$	99.72%	99.86%	99.16%
	$\mathcal{S}_f$	100%	99.58%	98.73%
	$\mathcal{S}_g$	99.72%	100%	99.15%
LLaMA-3.1-8B	$\mathcal{S}_{fg}$	100%	99.72%	99.72%
	$\mathcal{S}_f$	99.72%	98.03%	97.46%
	$\mathcal{S}_g$	100%	85.93%	99.72%
Qwen2.5-14B	$\mathcal{S}_{fg}$	99.72%	99.72%	99.44%
	$\mathcal{S}_f$	99.72%	99.44%	99.15%
	$\mathcal{S}_g$	99.86%	99.72%	99.86%
Yi-34B	$\mathcal{S}_{fg}$	99.86%	98.87%	99.15%
	$\mathcal{S}_f$	99.86%	99.86%	98.45%
	$\mathcal{S}_g$	99.72%	100%	99.58%

## H. Effect of Pretrained Encoder on Classifier Performance

To assess how classifier architecture affects unlearning trace detection, we evaluate a range of pretrained text encoders, following (BehnamGhader et al., 2024). We initialize classifiers with BERT (Devlin et al., 2019), T5 (Raffel et al., 2020), GPT-2 (Radford et al., 2019), and LLM2vec (BehnamGhader et al., 2024), each followed by a lightweight two-layer MLP head. Each classifier are trained to distinguish between responses generated by the original model and its -unlearned counterpart. As shown in **Tab. A6**, LLM2vec achieves the highest classification accuracy across all test scenarios. That’s why for adopting LLM2vec as the default classifier architecture.

To further probe how unlearning strength affects trace detectability across encoder architectures, we repeat our classification evaluation under the same mixed regime ( $\mathcal{S}_{fg}$ ) for both RMU- and NPO-unlearned Yi-34B outputs. **Tab. A6** and **Tab. A7** report accuracy when distinguishing original from unlearned responses using four different pretrained encoders. For RMU unlearning (Table A6), all encoders perform well on the WMDP “forget” data and MMLU “forget-irrelevant” data, but LLM2vec achieves the highest overall robustness, especially on UltraChat, where it attains 87.46% accuracy versus below 70% for the others. This validates our choice of LLM2vec as the default detector when unlearning traces are relatively subtle.

In stark contrast, Tab. A7 describes NPO unlearning yields near-perfect detection across both prompt types and all domains. Even the least robust encoder (T5) attains over 86% on UltraChat, while LLM2vec, GPT-2, and BERT all exceed 94% everywhere, with LLM2vec surpassing 99% on every test. This demonstrates that NPO’s more aggressive unlearning introduces globally visible artifacts, like raw token fragments and formatting anomalies, that make trace detection trivial, even on “forget-irrelevant” prompts where RMU traces often remain hidden.

Table A6. Classification accuracy for distinguishing original vs. RMU-unlearned responses using different pretrained sequence encoders. The source LLM is Yi-34B with RMU applied on the WMDP dataset. Settings mirror those in Tab. 2.

Classifier	WMDP	MMLU	UltraChat
LLM2vec	94.37%	95.77%	87.46%
T5	85.35%	82.96%	59.72%
GPT2	88.03%	96.06%	62.39%
BERT	88.59%	88.31%	69.15%

Table A7. Classification accuracy for distinguishing original vs. NPO-unlearned responses using different pretrained sequence encoders. Settings are consistent with Tab. A6.

Classifier	WMDP	MMLU	UltraChat
LLM2vec	99.86%	98.87%	99.15%
T5	99.29%	99.30%	86.20%
GPT2	99.72%	99.86%	96.90%
BERT	99.44%	99.58%	94.65%

## I. Multi-Class Classification Results

**Fig. A3** makes the NPO results visually striking. In panel (a), nearly every cell off the main diagonal is almost entirely dark, indicating that the classifier almost never confuses one class for another on WMDP “forget” prompts. The diagonal band itself is uniformly bright, reflecting the over 94% correct classification rates. Panel (b) shows a similarly unbroken diagonal on MMLU “forget-irrelevant” prompts, even Zephyr-7B, which under RMU had its signal buried in off-diagonals, now

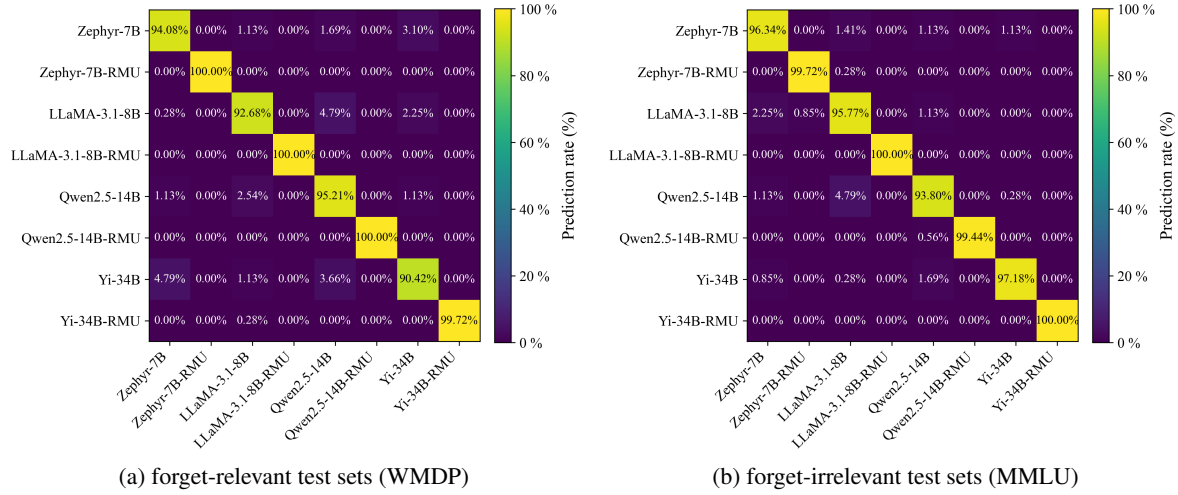


Figure A3. Confusion matrix for NPO-unlearning pair classification. Rows indicate true classes (original/NPO-unlearned model variants), and columns show predicted classes. Diagonal entries represent correct predictions; off-diagonals indicate misclassification rates under (a) WMDP and (b) MMLU test sets.

appears as a clean, high contrast stripe. The absence of any noticeable off-diagonal “leakage” in both panels underscores how aggressively NPO imprints a unique, domain-agnostic behavioral signature into every model variant, rendering multi-class identification effectively error-free.

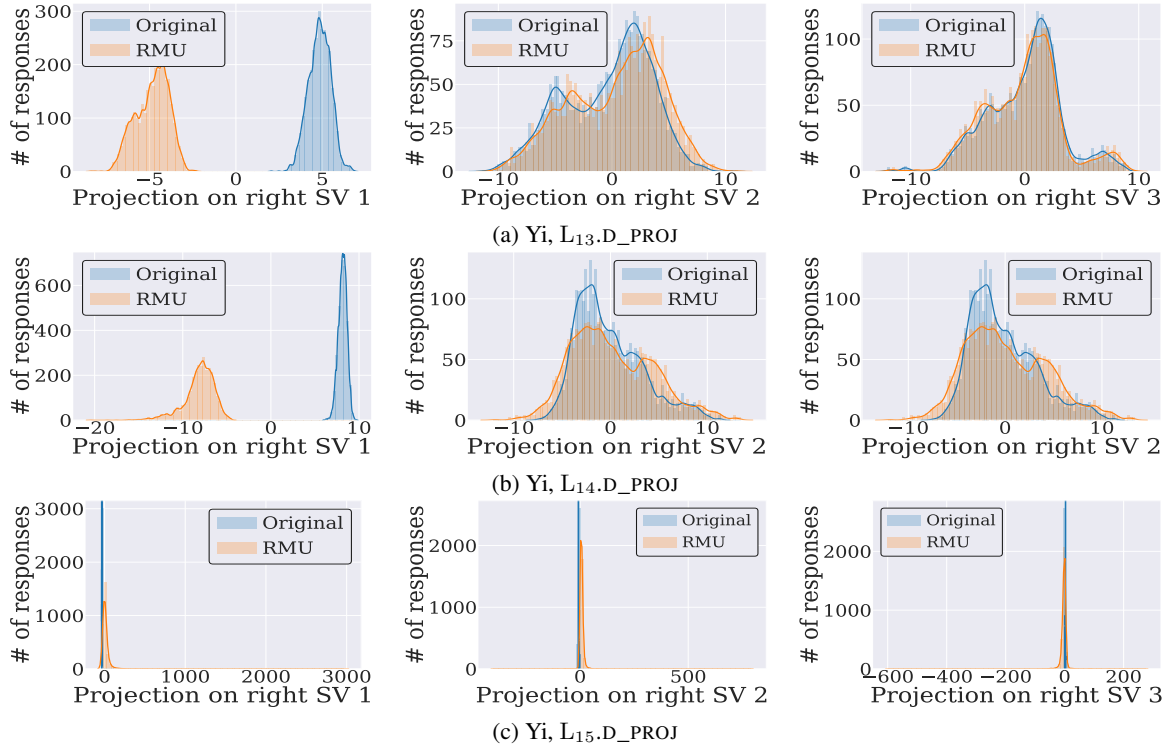


Figure A4. Projection of activations of Yi-34B-Chat from various layers for 3000 responses to MMLU onto the three leading right singular vectors for the original and unlearned model.  $L_i.D\_PROJ$  refers to activations extracted from the down-projection sublayer of the FFN in the  $i$ -th transformer block (a) are projections from layer 13, (b) are from layer 14, (c) are from layer 15.

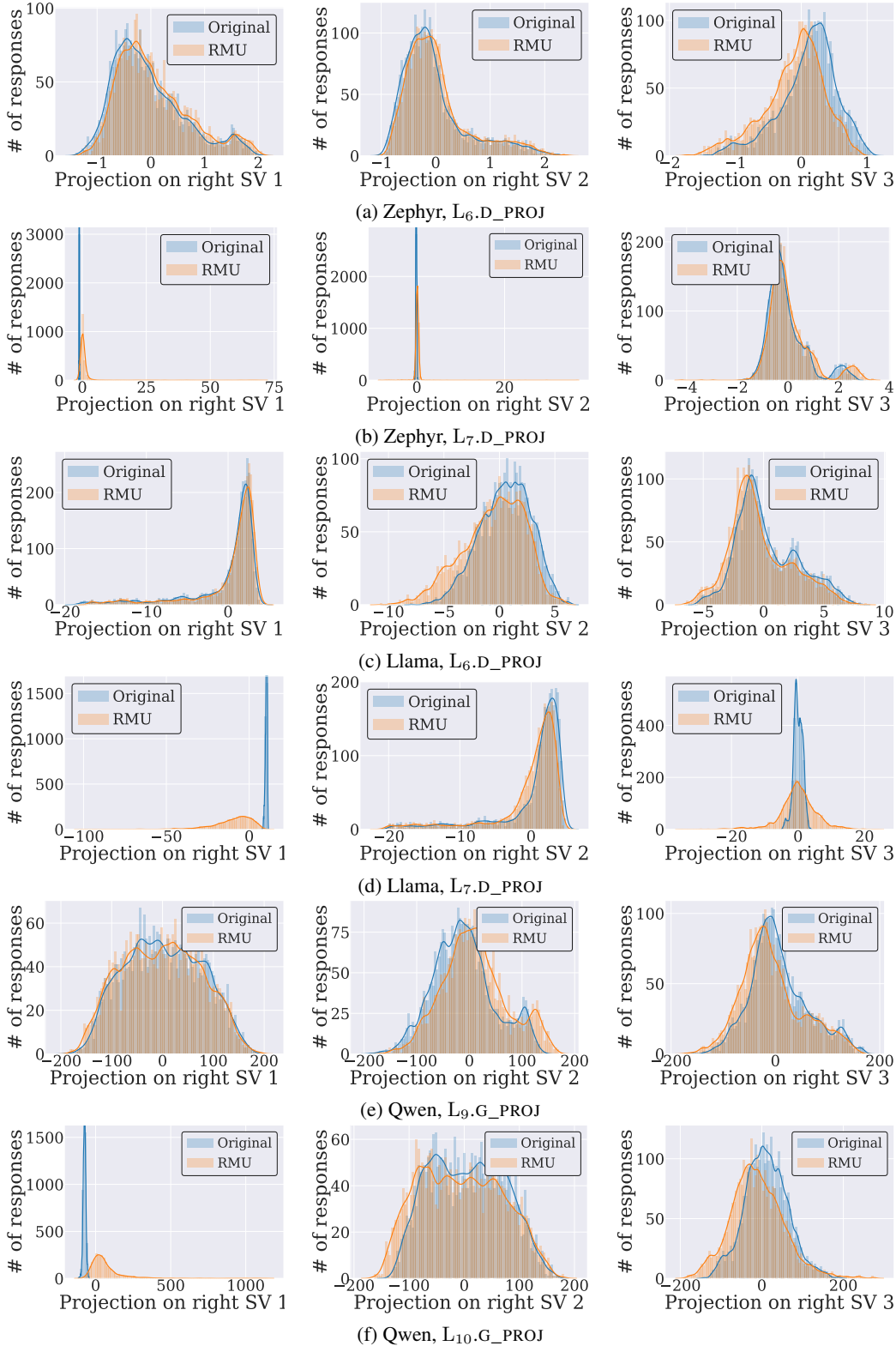


Figure A5. Projection of activations from various layers for 3000 responses to MMLU onto the three leading right singular vectors for the original and unlearned model.  $L_i.D\_PROJ$  refers to activations extracted from the down-projection sublayer of the FFN in the  $i$ -th transformer block, while  $L_i.G\_PROJ$  refers to activations extracted from the gate-projection sublayer of the FFN in the  $i$ -th transformer block (a,b) are projections for Zephyr-7B, (c,d) are for Llama3.1-8B, while (e,f) are for Qwen2.5-14B.