## 005 006 007 008 009 010 011 012 013 014 015 016

000 001

002 003 004

## **Gauss-Newton Unlearning for the LLM Era**

Anonymous Authors<sup>1</sup>

## Abstract

Large language models (LLMs) can learn to pro-011 duce sensitive outputs which model deployers may 012 wish to reduce, motivating the use of output suppression (LLM unlearning) methods. We demonstrate that taking only a few uphill Gauss-Newton 015 steps on a forget set provides a conceptually simple, state-of-the-art unlearning algorithm that is underexplored in the LLM literature. We show 018 that these steps can be efficiently and accurately 019 implemented for LLMs using parametric Hessian 020 approximations such as K-FAC. We call this approach K-FAC for Distribution Erasure (K-FADE). Our evaluations demonstrate that K-FADE performs competitively with or better than previous unlearning approaches for LLMs across standard 025 benchmarks. Specifically, K-FADE approximates the output distribution of models re-finetuned with 027 certain data excluded on the ToFU unlearning 028 benchmark. K-FADE also effectively suppresses 029 outputs from a specific distribution while mini-030 mally altering the model's outputs on non-targeted data from the WMDP benchmark.

### Introduction

034

035

049

050

051

052

053

054

Large Language Models (LLMs) pre-train on large swaths of the internet, learning subdistributions that a model deployer 038 may not desire. For example, models can memorize sensitive 039 information (Huang et al., 2022; Carlini et al., 2021), like emails and phone numbers, or produce content that may be 041 useful in the construction of chemical, biological, radiological, and nuclear (CBRN) weapons (Li et al., 2024). Output 043 suppression techniques (LLM unlearning) attempt to mitigate these and other socio-technical harms by decreasing the prob-045 ability of producing outputs from certain subdistributions 046 without the need to retrain the model from scratch (Liu et al., 047 2024b; Cooper et al., 2024).

Similar to past LLM unlearning papers, we aim to satisfy two desiderata (Li et al., 2024; Zhang et al., 2024; Fan et al., 2024; Yao et al., 2024b). First, we want the model to perform poorly on a *forget set* that measures the model's ability to produce unwanted information, a desideratum we call output suppression (Cooper et al., 2024; Liu et al., 2024b). Second, we want the model's behavior to remain as close as possible to the initial model in other settings; we call this desideratum specificity. We refer to the data we use to estimate and evaluate this specificity constraint as the retain set. Past methods struggle to optimize for output suppression while maintaining specificity. They typically rely on many steps using explicit noisy KL loss penalties (Maini et al., 2024; Yao et al., 2024b;a; Li et al., 2024; Gandikota et al., 2024) or low fidelity second-order methods (Jia et al., 2024) to maintain the specificity constraint. These limitations often cause significant changes in model outputs well beyond the intended forget distribution. We provide a more extensive discussion of related work in Appendix B.

Building on advances in parametric Hessian estimation (Martens & Grosse, 2015; George et al., 2018; Grosse et al., 2023), we show how Gauss-Newton updates, a conceptually simple traditional ML unlearning method, scale to LLMs with billions of parameters. Researchers originally proposed Gauss-Newton ascent steps as an unlearning method for linear models (Guo et al., 2019; Izzo et al., 2021). However, by linearizing a neural network around its final parameters (Jia et al., 2023), we can apply an analogous technique to large neural networks, which we observe is equivalent to natural gradient ascent. Despite this technique's conceptual simplicity, the difficulty in approximating the necessary inverse Hessian vector products has limited its application to large neural networks like LLMs. Progress in accurately estimating the Gauss-Newton Hessian, e.g., K-FAC (Martens & Grosse, 2015) and EK-FAC (George et al., 2018), and work on scaling these techniques to LLMs (Grosse et al., 2023), now makes this approach feasible; we find these sophisticated estimators are necessary for good performance on unlearning benchmarks (see ablations in Appendix E). We call our approach K-FAC for Distribution Erasure (K-FADE).

We validate our method on standard unlearning benchmarks. First, we evaluate it on the Weapons of Mass Destruction Proxy (WMDP) benchmark (Li et al., 2024), which measures

<sup>&</sup>lt;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.

055 a method's ability to suppress proxies for "hazardous" content while maintaining broad knowledge and fluency. After 057 unlearning on WMDP, we show that K-FADE matches state-058 of-the-art results in knowledge suppression while preserving 059 "benign" knowledge (Hendrycks et al., 2021), achieving supe-060 rior fluency (Zheng et al., 2023) and superior specificity, as 061 measured by a novel metric based on the model's increase in 062 KL divergence on diverse instruction following data (Taori 063 et al., 2023). Second, we evaluate it on the Test of Fictitious 064 Unlearning (ToFU) (Maini et al., 2024), which benchmarks 065 the ability to remove sensitive information about individuals 066 while preserving nonsensitive information. We find that a 067 single, correctly implemented Gauss-Newton step delivers 068 a Pareto improvement in forget quality and model utility, 069 representing a new state-of-the-art on the ToFU benchmark 070 that has relatively few hyperparameters. 071

## Background: Gauss-Newton Unlearning

074 Sometimes, we know precisely which data caused a gen-075 erative model to produce particular sensitive outputs (e.g., 076 obscure facts about individuals or memorized text (Carlini 077 et al., 2021)). In such cases, the forget set is clear, and ap-078 proximating the output distribution of a model never trained 079 on this data (approximate unlearning) serves as an appropriate gold standard. In many cases, however, there is no 081 clear forget set. In such cases our primary objective is output 082 suppression where we aim to suppress particular subdistribu-083 tions while minimally changing model behavior outside this 084 synthetic forget distribution.

A common approach to output suppression combines decreasing the probability of outputs from a forget set  $D_f$ , using a forgetting loss  $L_F(\theta)$ , while maintaining performance on a retain set  $D_r$  using a KL penalty (Kurmanji et al., 2023; Gandikota et al., 2024; Maini et al., 2024; Liu et al., 2022):

085

086

087

089

090

091

092

093

094

$$\mathcal{L}_F(\theta) + \gamma \mathbb{E}_{x \sim D_r} \left[ \mathrm{KL}(p(.|f(x;\theta^*)), p(.|f(x;\theta))) \right]$$
(1)

095 Past methods in the LLM unlearning literature have focused 096 on optimizing this objective largely using first-order methods 097 (Zhang et al., 2024; Fan et al., 2024; Maini et al., 2024). 098 In our work, we take advantage of the fact that optimizing 099 this objective is locally equivalent to optimizing the natu-100 ral gradient of the forget loss  $L_F(\theta)$ . The natural gradient direction is given by  $G_{\theta}^{-1} \nabla_{\theta} \mathcal{L}(\theta)$ , where  $G_{\theta}$  is the Fisher Information Matrix (FIM). For all the models we consider in this paper the FIM is equivalent to the Gauss-Newton 104 Hessian (GNH) meaning that this update is equivalent to a 105 single Gauss-Newton step. Thus we call this general strategy 106 Gauss-Newton Unlearning. The advantage of this second or-107 der step over taking many first-order steps is that it allows us to consider the effects of the update on a large retain set  $D_r$ 109

at every step leading to stable and highly targeted unlearning updates.

Interestingly, taking Gauss-Newton steps has an interpretation as an approximate unlearning method for linear models that minimize a strictly convex loss (Guo et al., 2019; Izzo et al., 2021). We can apply an analogous update to neural networks by linearizing them around their final parameters, i.e., approximating the network using the Jacobian of its outputs with respect to weights (Bae et al., 2022; Jia et al., 2023). Under this linearization, the unlearning update direction is again proportional to the natural gradient with respect to the training loss on the forget set. We provide a more detailed discussion of how approximate retraining, output suppression, and the natural gradient relate in Appendix C.

### Methods

In this section we explore how we can apply parametric Hessian approximations to perform the needed inverse Hessian vector products (iHVPs) for Gauss-Newton unlearning.

Efficient second-order approximations. Explicitly representing the entire Hessian matrix is impractical for LLMs. While there are techniques that seek to approximate iHVPs without constructing the Hessian (e.g. Martens (2010)), these typically don't scale to large datasets and models. Since working with LLMs and large retain sets, is our explicit goal, we focus on parametric approximations. Diagonal approximations (Becker, 1988; Liu et al., 2024a) are cheap but miss parameter inter-dependencies and we find that they are not effective for unlearning in ToFU (see Appendix E). Thus we turn to methods like Kronecker-Factored Approximate Curvature (K-FAC) (Martens & Grosse, 2015) and eigenvaluecorrected K-FAC (EK-FAC) (George et al., 2018; Grosse et al., 2023) which provide compact approximations. We use the same strategies to handle weight sharing and K-FAC/EK-FAC factor fits as Grosse et al. (2023) and we build our implementation on the excellent CurvLinOps library (Dangel et al., 2025). All of these parametric approximations require the use of *damping* as many of the GNH eigenvalues are very small. This means that in practice we invert  $G_{\theta} + \lambda I$  instead of  $G_{\theta}$  where  $\lambda$  is the damping parameter.

Beyond using Gauss-Newton steps to do natural gradient ascent on the forget set, there are several additional tricks that make K-FADE unlearning effective in different settings, as we describe below.

**Suppression objective.** We consider two objective functions for  $\mathcal{L}_F$  to achieve output suppression: increasing the margin and increasing the cross entropy. When approximating retraining in linear models, the forgeting objective is simpling the increasing the training objective (Guo et al., 2019). And indeed in tasks where matching retraining is desirable, the cross entropy is effective. However, we find that the margin

Gauss-Newton Unlearning for the LLM Era

Model	Method	WMDP		Model Utility			
11100001	method	Bio ↓	Cyber $\downarrow$	MMLU↑	MT-Bench↑	$D_{KL} \times 10^{-2} \downarrow$	
	Original	$64.3^{\pm 1.3}$	<b>44.7</b> <sup><math>\pm 1.0</math></sup>	$58.4^{\pm0.4}$	7.2	0	
Zephr-7b- $\beta$	ELM RMU K-FADE (Ours)	$\begin{array}{c} \textbf{29.8}^{\pm 1.3} \\ \underline{30.4}^{\pm 1.3} \\ \overline{30.1}^{\pm 1.3} \end{array}$	$\frac{\underline{27.3}^{\pm 1.0}}{27.1^{\pm 1.0}}$ $27.7^{\pm 1.0}$	56.7 <sup>±0.4</sup> 57.5 <sup>±0.4</sup> 57.2 <sup>±0.4</sup>	$\frac{6.86}{6.71^{\pm 0.03}}$ 6.71 <sup>±0.07</sup> 6.91 <sup>±0.04</sup>	<b>6.7</b> [6.3–6.9] <b>5.3</b> [4.4–6.1] <b>2.9</b> [2.4–3.5]	

164

119

120

Table 1. K-FADE attains state-of-the-art performance in repressing hazardous knowledge while having better unlearning specificity. Like ELM (Gandikota et al., 2024) and RMU (Li et al., 2024), K-FADE reduces model performance on WMDP's Bio and Cyber significantly, while retaining performance on MMLU (Hendrycks et al., 2021) and MT-Bench (Zheng et al., 2023). However, we see K-FADE preserves model behavior on a diverse instruction following dataset, alpaca (Taori et al., 2023), as measured by KL divergence, significantly better than the baselines. For multiple choice questions and MT-Bench (n=5) we report stderr, on the mean KL  $D_{KL}$  over alpaca we report 95% bootstrapped CIs. Methods that are not significantly different from the best are underlined.

is generally more effective for tasks requiring high unlearning specificity and multiple unlearning steps. When we refer to the per example margin (Park et al., 2023), it is defined as,  $\ell^{(\text{margin})}(z; y) = z_y - \log \sum_{i \neq y} \exp(z_i)$ . We sample batches of data without replacement from the forget set to compute our objective gradient.

**Step size.** Tuning the step size  $\alpha$  is essential for effective unlearning. Our motivating example of convex unlearning suggests simply adding the inverse Hessian vector product  $r:= \left(\tilde{G}_{\theta}\right)^{-1}g$  to the model parameters  $\theta' \leftarrow \theta + r$  at each step, where  $g := \nabla_{\theta} \mathcal{L}_F$ . In practice, we find that the natural gradient offers a better heuristic for step size selection. 138 Borrowing from the geometric interpretation of the natural 139 gradient (Martens, 2020), we ensure that each step causes ap-140 proximately constant KL divergence by ensuring the step has 141 constant norm  $\alpha$  under the  $G_{\theta}$  inner product. This method is 142 very similar to the technique used in Ba et al. (2017). We find 143 that this approach makes unlearning over multiple steps more 144 stable and decouples the effects of changing the damping 145 parameter  $\lambda$  and step size  $\alpha$ .

What to fit the Hessian on. Implementations on linear mod-147 els, suggest that we should only be fitting the Gauss-Newton 148 Hessian using only the retain set (Guo et al., 2019). Support-149 ing this, we find that including the forget set in the Hessian 150 computation generally reduces the specificity of the unlearn-151 ing updates. For our experiments, we fit the Hessian on only 152 the retain set. See our ablation study (Appendix E) for the 153 effects of including the forget set in the GNH computation. 154

155 Components targeted. Following Grosse et al. (2023), our 156 method only targets the weights in the affine transformations 157 in the model's feed-forward layers. This still encompasses a 158 large fraction of the model's total parameters (e.g. 78% for 159 Mistral-7b (Jiang et al., 2023)). For some experiments, we 160 target only a subset of these layers (e.g. WMDP (Li et al., 161 2024)) to improve unlearning specificity. In Appendix D we 162 describe when and how we implement this targeting. 163

### **Experiments**

In our experiments, we aim to address how K-FADE compares to baselines for output suppression while maintaining specificity (Li et al., 2024) and matching retraining (Maini et al., 2024). We include additional experiments exploring Hessian ablations and fine-tuning attacks in Appendix E.

#### RQ1: CAN K-FADE SUPPRESS HARMFUL KNOWLEDGE?

K-FADE provides state-of-the-art output suppression (Figure 1) while providing better specificity.

**WMDP.** The Weapons of Mass Destruction Proxy (WMDP) benchmark (Li et al., 2024) assesses a model's ability to output proxies for hazardous knowledge in cybersecurity (WMDP Cyber), bio-weapons (WMDP Bio), and chemical weapons (WMDP Chem). The benchmark uses multiple-choice questions and provides forget sets of relevant documents for each domain. We use Wikitext (Merity et al., 2017) as our retain set, following Li et al. (2024). The specificity of unlearning is measured using MMLU (Hendrycks et al., 2021) which measures general knowledge, MT-Bench (Zheng et al., 2023) which measures "fluency" as judged by gpt-4.

**Specificity Evaluation.** We introduce an additional specificity evaluation where we measure the KL divergence from the base model to the unlearned models on 30000 instructions from the Alpaca dataset (Taori et al., 2023), generating completions from zephyr-7b- $\beta$  (Tunstall et al., 2023). We report the average KL per-token only on the completions, not the instructions. Unlike MT-Bench (Zheng et al., 2023), it does not depend on an additional LLM as an auto-grader. Generally, we find that observing which completions have high KL is useful for understanding the side effects of unlearning methods, e.g., K-FADE models unlearned on WMDP Bio will not discuss experiments involving mice, and RMU models generally refuse to discuss COVID-19.

Baselines. We compare K-FADE to RMU (Li et al., 2024)



*Figure 1.* **One-step of K-FADE outperforms the state of the art in unlearning on the TOFU dataset.** The left figure shows performance when unlearning 5% of the authors bios, the right figure displays the results for unlearning 10%. Forget quality measures how close the distribution of model responses about the unlearned authors is to a model that was never trained on these authors. The model utility is then the model's ability to recall facts about fictitious authors from the retain set and real authors. K-FADE effectively outperforms both of the baseline methods provided in the original TOFU paper (Maini et al., 2024) and a recent state of the art method simNPO (Fan et al., 2024).

and ELM (Gandikota et al., 2024). RMU disrupts activations
relevant to the forget set while minimizing L2 distance in
activation space to preserve performance. ELM combines
a steering loss inspired by classifier-free guidance with KL
divergence and fluency penalties. We directly compare to the
model checkpoints provided by the authors of the RMU and
ELM papers.

K-FADE achieves state-of-the-art specificity. K-FADE 189 achieves strong output suppression, matching RMU and ELM 190 on WMDP-Bio and WMDP-Cyber (see Table 1). In terms of 191 specificity, performance on MMLU (Hendrycks et al., 2021) is similar to RMU (Li et al., 2024), and better than ELM 193 (Gandikota et al., 2024) indicating that most knowledge is preserved. In terms of fluency as measured by MT-Bench 195 196 (Zheng et al., 2023), K-FADE is significantly better than both ELM and RMU. Additionally, the average KL divergence 197 between a model unlearned with K-FADE and the base model (zephyr-7b- $\beta$  (Tunstall et al., 2023)) is 40% lower than the 199 next best method RMU. Interestingly, ELM, RMU, and K-200 FADE show distinct distributional effects: ELM changes the output distribution over nearly all documents while not having a long tail of increased KL divergence; RMU shows 203 a more targeted effect with a distinct long tail of radically 204 changed completions; K-FADE behaves similarly to RMU 206 in the tail but shows lower KL divergence in the head of the distribution (Figure 2).

#### 209 RQ2: CAN K-FADE APPROXIMATE RETRAINING?

K-FADE approximately removes the effects of fine-tuning datapoints, as demonstrated by ToFU using a single Gauss-Newton step. We give experiment details in Appendix D.2.

reference models) and model utility (ability to recall world knowledge, answer about real authors, and retained fictitious authors). We experiment with unlearning 5% and 10% of these Q&A pairs.

**Baselines.** We compare to baselines from the ToFU paper: **Grad. Ascent, Grad. Diff.**, and **KL min.** As well as strong recent baselines **simNPO** (Fan et al., 2024) with default hyperparameters ( $\beta = 2.5$ , NPO coefficient=0.1375 for 5%;  $\beta = 4.5$ , NPO coefficient=0.125 for 10%) and **SOUL** (Jia et al., 2024) where we again use their default hyper-parameters which they only provide for the 10% set <sup>1</sup>

**One Gauss-Newton step achieves state-of-the-art on ToFU.** Our method achieves state-of-the-art forget quality on the challenging 10% forget set, outperforming the original ToFU baselines, simNPO (Fan et al., 2024). We do this while achieving comparable model utility on both the 5% and 10% sets (Figure 1, Table 2).

#### Conclusions

We have shown that a conceptually simple unlearning algorithm, Gauss-Newton ascent steps, can be efficiently scaled to LLMs, is state of the art on multiple benchmarks, and is particularly effective at preserving the model's performance on non-targeted data. On one benchmark, we even find that a single Gauss-Newton step can outperform the previous state of the art, allowing for Hessian caching that significantly reduces the cost of hyperparameter tuning. Finally, we introduced a novel measure of LLM unlearning specificity: evaluating the KL divergence between the base and unlearned model's outputs on thousands of benign completions.

208

175

176

177

178

<sup>ToFU. The ToFU benchmark (Maini et al., 2024) contains
questions and answers about fictitious authors, with models
finetuned on these Q&A pairs. The goal is to unlearn facts
about a subset of authors. Performance measures include
forget quality (similarity between unlearned and retrained</sup> 

<sup>&</sup>lt;sup>1</sup>SOUL (Jia et al., 2024) gets a much lower forget quality than originally reported in their paper in our results. This is because we and the original ToFU (Maini et al., 2024) paper use a messure of forget quality better aligned with matching the output distribution of re-fine-tuned models.

### 220 References

221

222

223

258

259

261

272

273

274

- Amari, S.-I. Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276, 1998.
- Ba, J., Grosse, R., and Martens, J. Distributed second-order
  optimization using kronecker-factored approximations. In *International Conference on Learning Representations*,
  2017. URL https://openreview.net/forum?
  id=SkkTMpjex.
- Bae, J., Ng, N. H., Lo, A., Ghassemi, M., and Grosse, R. B.
  If influence functions are the answer, then what is the
  question? In Oh, A. H., Agarwal, A., Belgrave, D., and
  Cho, K. (eds.), Advances in Neural Information Processing
  Systems, 2022. URL https://openreview.net/
  forum?id=hzbguA9zMJ.
- Becker, S. Improving the convergence of backpropagation learning with second order method. In *Proceedings of the 1988 Connectionist Models Summer School, San Mateo, CA.* Morgan Kaufmann, 1988.
- Carlini, N., Tramer, F., Wallace, E., Jagielski, M., Herbert-Voss, A., Lee, K., Roberts, A., Brown, T., Song, D., Erlingsson, U., et al. Extracting training data from large language models. In *30th USENIX Security Symposium* (USENIX Security 21), pp. 2633–2650, 2021.

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.

- Dangel, F., Eschenhagen, R., Ormaniec, W., Fernandez, A.,
  Tatzel, L., and Kristiadi, A. Position: Curvature matrices
  should be democratized via linear operators, 2025. URL
  https://arxiv.org/abs/2501.19183.
  - Deeb, A. and Roger, F. Do unlearning methods remove information from language model weights?, 2024. URL https://arxiv.org/abs/2410.08827.
- Fan, C., Liu, J., Lin, L., Jia, J., Zhang, R., Mei, S., and
  Liu, S. Simplicity prevails: Rethinking negative preference optimization for LLM unlearning. In *Neurips Safe Generative AI Workshop 2024*, 2024. URL https:
  //openreview.net/forum?id=pVACX02m0p.
- Gandikota, R., Feucht, S., Marks, S., and Bau, D. Erasing conceptual knowledge from language models. *CoRR*, abs/2410.02760, 2024. URL https://doi.org/10.48550/arXiv.2410.02760.
  - George, T., Laurent, C., Bouthillier, X., Ballas, N., and Vincent, P. Fast approximate natural gradient descent in a

kronecker factored eigenbasis. Advances in neural information processing systems, 31, 2018.

- Golatkar, A., Achille, A., and Soatto, S. Eternal Sunshine of the Spotless Net: Selective Forgetting in Deep Networks, March 2020. URL http://arxiv.org/abs/1911. 04933. arXiv:1911.04933 [cs].
- Grosse, R., Bae, J., Anil, C., Elhage, N., Tamkin, A., Tajdini, A., Steiner, B., Li, D., Durmus, E., Perez, E., Hubinger, E., Lukošiūtė, K., Nguyen, K., Joseph, N., McCandlish, S., Kaplan, J., and Bowman, S. R. Studying large language model generalization with influence functions, 2023. URL https://arxiv.org/abs/2308.03296.
- Gu, K., Rashid, M. R. U., Sultana, N., and Mehnaz, S. Second-Order Information Matters: Revisiting Machine Unlearning for Large Language Models, March 2024. URL http://arxiv.org/abs/2403. 10557. arXiv:2403.10557 [cs].
- Guo, C., Goldstein, T., Hannun, A., and Van Der Maaten, L. Certified data removal from machine learning models. arXiv preprint arXiv:1911.03030, 2019.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021. URL https: //openreview.net/forum?id=d7KBjmI3GmQ.
- Huang, J., Shao, H., and Chang, K. C.-C. Are large pretrained language models leaking your personal information? In Goldberg, Y., Kozareva, Z., and Zhang, Y. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 2038–2047, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp. 148. URL https://aclanthology.org/2022. findings-emnlp.148/.
- Ilharco, G., Ribeiro, M. T., Wortsman, M., Gururangan, S., Schmidt, L., Hajishirzi, H., and Farhadi, A. Editing models with task arithmetic. In *Proceedings of the 11th International Conference on Learning Representations*, 2023. URL https://arxiv.org/abs/2212.04089.
- Izzo, Z., Smart, M. A., Chaudhuri, K., and Zou, J. Approximate data deletion from machine learning models. In *International Conference on Artificial Intelligence and Statistics*, pp. 2008–2016. PMLR, 2021.
- 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.

- 275 Jia, J., Liu, J., Ram, P., Yao, Y., Liu, G., Liu, Y., Sharma, P., Liu, H., Li, Z., Hall, D., Liang, P., and Ma, T. Sophia: A 276 and Liu, S. Model sparsity can simplify machine unlearn-277 ing. Advances in Neural Information Processing Systems, 278 36:51584-51605, 2023.
- 279 Jia, J., Zhang, Y., Zhang, Y., Liu, J., Runwal, B., Diffend-280 erfer, J., Kailkhura, B., and Liu, S. SOUL: Unlocking 281 the Power of Second-Order Optimization for LLM Un-282 learning, June 2024. URL http://arxiv.org/abs/ 283 2404.18239. arXiv:2404.18239 [cs]. 284
- 285 Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., 286 Chaplot, D. S., de las Casas, D., Bressand, F., Lengyel, 287 G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, 289 T., and Sayed, W. E. Mistral 7b, 2023. URL https: 290 //arxiv.org/abs/2310.06825. 291
- Kingma, D. P. and Ba, J. Adam: A method for stochastic 292 optimization. In Proceedings of the 3rd International Con-293 ference on Learning Representations, ICLR, 2015. URL 294 http://arxiv.org/abs/1412.6980. 295
- 296 Koh, P. W. and Liang, P. Understanding Black-297 box Predictions via Influence Functions, December 2017. URL http://arxiv.org/abs/1703. 299 04730. arXiv:1703.04730 [cs, stat]. 300
- Kurmanji, M., Triantafillou, P., Hayes, J., and Triantafillou, 301 E. Towards unbounded machine unlearning. In Thirty-302 seventh Conference on Neural Information Processing 303 Systems, 2023. URL https://openreview.net/ 304 forum?id=OveBaTtUAT. 305
- 306 Li, N., Pan, A., Gopal, A., Yue, S., Berrios, D., Gatti, A., Li, 307 J. D., Dombrowski, A.-K., Goel, S., Phan, L., Mukobi, G., 308 Helm-Burger, N., Lababidi, R., Justen, L., Liu, A. B., 309 Chen, M., Barrass, I., Zhang, O., Zhu, X., Tamirisa, R., Bharathi, B., Khoja, A., Zhao, Z., Herbert-Voss, A., 311 Breuer, C. B., Marks, S., Patel, O., Zou, A., Mazeika, 312 M., Wang, Z., Oswal, P., Lin, W., Hunt, A. A., Tienken-313 Harder, J., Shih, K. Y., Talley, K., Guan, J., Kaplan, R., 314 Steneker, I., Campbell, D., Jokubaitis, B., Levinson, A., 315 Wang, J., Qian, W., Karmakar, K. K., Basart, S., Fitz, S., 316 Levine, M., Kumaraguru, P., Tupakula, U., Varadharajan, 317 V., Wang, R., Shoshitaishvili, Y., Ba, J., Esvelt, K. M., 318 Wang, A., and Hendrycks, D. The WMDP Benchmark: 319 Measuring and Reducing Malicious Use With Unlearning, 320 May 2024. URL http://arxiv.org/abs/2403. 321 03218. arXiv:2403.03218 [cs].
- Li, Y., Bubeck, S., Eldan, R., Del Giorno, A., Gunasekar, 323 S., and Lee, Y. T. Textbooks are all you need ii: phi-1.5 324 technical report. arXiv preprint arXiv:2309.05463, 2023. 325
- Liu, B., Liu, Q., and Stone, P. Continual learning and private 327 unlearning. In Conference on Lifelong Learning Agents, 328 pp. 243-254. PMLR, 2022. 329

- scalable stochastic second-order optimizer for language model pre-training, 2024a. URL https://arxiv. org/abs/2305.14342.
- Liu, S., Yao, Y., Jia, J., Casper, S., Baracaldo, N., Hase, P., Yao, Y., Liu, C. Y., Xu, X., Li, H., Varshney, K. R., Bansal, M., Koyejo, S., and Liu, Y. Rethinking Machine Unlearning for Large Language Models, July 2024b. URL http://arxiv.org/abs/2402. 08787. arXiv:2402.08787 [cs].
- Maini, P., Feng, Z., Schwarzschild, A., Lipton, Z. C., and Kolter, J. Z. TOFU: A Task of Fictitious Unlearning for LLMs, January 2024. URL http://arxiv.org/ abs/2401.06121. arXiv:2401.06121 [cs].
- Martens, J. Deep learning via hessian-free optimization. In Proceedings of the 27th International Conference on International Conference on Machine Learning, ICML'10, pp. 735-742, Madison, WI, USA, 2010. Omnipress. ISBN 9781605589077.
- Martens, J. New insights and perspectives on the natural gradient method. Journal of Machine Learning Research, 21(146):1-76, 2020. URL http://jmlr.org/ papers/v21/17-678.html.
- Martens, J. and Grosse, R. Optimizing neural networks with kronecker-factored approximate curvature. In Bach, F. and Blei, D. (eds.), Proceedings of the 32nd International Conference on Machine Learning, volume 37 of Proceedings of Machine Learning Research, pp. 2408-2417, Lille, France, 07-09 Jul 2015. PMLR. URL https://proceedings.mlr.press/v37/ martens15.html.
- Merity, S., Xiong, C., Bradbury, J., and Socher, R. Pointer sentinel mixture models. In International Conference on Learning Representations, 2017. URL https:// openreview.net/forum?id=Byj72udxe.
- Park, S. M., Georgiev, K., Ilyas, A., Leclerc, G., and Mądry, A. Trak: attributing model behavior at scale. In Proceedings of the 40th International Conference on Machine Learning, pp. 27074-27113, 2023.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E. Z., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. Pytorch: An imperative style, high-performance deep learning library. CoRR, abs/1912.01703, 2019. URL http://arxiv.org/ abs/1912.01703.

330 331 332 333	Qi, X., Wei, B., Carlini, N., Huang, Y., Xie, T., He, L., Jagielski, M., Nasr, M., Mittal, P., and Henderson, P. On evaluating the durability of safeguards for open-weight llms. <i>arXiv preprint arXiv:2412.07097</i> , 2024.	Zhang, R., Lin, L., Bai, Y., and Mei, S. Negative prefer- ence optimization: From catastrophic collapse to effective unlearning, 2024. URL https://arxiv.org/abs/ 2404.05868.
334 335 336 337 338 339 340 341 342	Rosati, D., Wehner, J., Williams, K., Bartoszcze, L., Gon- zales, R., carsten maple, Majumdar, S., Sajjad, H., and Rudzicz, F. Representation noising: A defence mechanism against harmful finetuning. In <i>The Thirty-eighth Annual</i> <i>Conference on Neural Information Processing Systems</i> , 2024. URL https://openreview.net/forum? id=eP9auEJqFg.	Zheng, L., Chiang, WL., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li13, D., Xing35, E. P., et al. Judging llm-as-a-judge with mt-bench and chatbot arena. arXiv preprint arXiv:2306.05685, 2023.
343 344 345 346 347 348	Shi, W., Lee, J., Huang, Y., Malladi, S., Zhao, J., Holtzman, A., Liu, D., Zettlemoyer, L., Smith, N. A., and Zhang, C. MUSE: Machine Unlearning Six-Way Evaluation for Language Models, July 2024. URL http://arxiv. org/abs/2407.06460. arXiv:2407.06460 [cs].	
349 350 351 352	Singh, S. P. and Alistarh, D. Woodfisher: Efficient second- order approximation for neural network compression. Ad- vances in Neural Information Processing Systems, 33: 18098–18109, 2020.	
353 354 355 356 357 358	<ul> <li>Tamirisa, R., Bharathi, B., Phan, L., Zhou, A., Gatti, A., Suresh, T., Lin, M., Wang, J., Wang, R., Arel, R., Zou, A., Song, D., Li, B., Hendrycks, D., and Mazeika, M. Tamper-resistant safeguards for open-weight llms, 2024. URL https://arxiv.org/abs/2408.00761.</li> </ul>	
359 360 361 362 363 364	Taori, R., Gulrajani, I., Zhang, T., Dubois, Y., Li, X., Guestrin, C., Liang, P., and Hashimoto, T. B. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/ stanford_alpaca, 2023.	
365 366 367 368 369	Thudi, A., Jia, H., Shumailov, I., and Papernot, N. On the necessity of auditable algorithmic definitions for machine unlearning. In <i>31st USENIX security symposium (USENIX Security 22)</i> , pp. 4007–4022, 2022.	
370 371 372 373 374 375	Tunstall, L., Beeching, E., Lambert, N., Rajani, N., Rasul, K., Belkada, Y., Huang, S., von Werra, L., Fourrier, C., Habib, N., Sarrazin, N., Sanseviero, O., Rush, A. M., and Wolf, T. Zephyr: Direct distillation of lm alignment, 2023. URL https://arxiv.org/abs/2310.16944.	
376 377 378 379 380	Yao, J., Chien, E., Du, M., Niu, X., Wang, T., Cheng, Z., and Yue, X. Machine Unlearning of Pre-trained Large Language Models, May 2024a. URL http://arxiv. org/abs/2402.15159. arXiv:2402.15159 [cs].	
381 382 383 384	Yao, Y., Xu, X., and Liu, Y. Large Language Model Un- learning, February 2024b. URL http://arxiv.org/ abs/2310.10683. arXiv:2310.10683 [cs].	
		7

## 385 Appendix

# A. Additional visualizations of the results

		Summary		ROUGE			
		Forget Q. ↑	Utility ↑	Retain ↑	Auth ↑	World ↑	Forget $\downarrow$
Split	Method						-
Forget 5%	Retain Model	1.00	0.60	0.98	0.95	0.89	0.40
	Grad. Ascent	4.61e-07	0.44	0.36	0.82	0.90	0.31
	Grad. Diff.	1.18e-02	0.40	0.26	0.57	0.82	0.13
	KL Min.	1.46e-07	0.43	0.35	0.80	0.90	0.30
	SimNPO	0.63	0.51	0.44	0.87	0.88	0.33
	K-FADE (Ours)	0.87	0.57	0.61	0.91	0.85	0.31
Forget 10%	Retain Model	1.00	0.61	0.98	0.92	0.90	0.41
	Grad. Ascent	2.19e-16	0.63	0.70	0.94	0.92	0.59
	Grad. Diff.	3.34e-04	0.07	0.09	0.17	0.65	0.07
	KL Min.	1.06e-16	0.63	0.72	0.94	0.91	0.61
	SimNPO	2.08e-02	0.52	0.44	0.85	0.86	0.37
	SOUL Grad. Diff.	5.56e-14	0.58	0.45	0.60	0.86	0.02
	K-FADE (Ours)	0.85	0.57	0.52	0.90	0.86	0.32

*Table 2.* **TOFU Results One gauss-newton step using K-FADE is a pareto-improvement over the baseline methods across model
<b>utility and forget quality.** Here we show the checkpoints that have the largest product of Forget Quality and Model Utility from Figure
1. K-FADE achieves a state of the art forget quality, which attempts to approximate the similarity of the model's outputs on the forget
distribution to a model only trained on the retain set (Maini et al., 2024). It does this while preserving model utility on the challenging
10% forget set.

Gauss-Newton Unlearning for the LLM Era



*Figure 2.* **K-FADE perturbs the model's behavior on unrelated data less than strong baselines** like ELM (Gandikota et al., 2024) and RMU (Li et al., 2024). The plot shows the distribution of KL divergences, on completions generated from prompts in the alpaca dataset, between zephyr-7b- $\beta$  and models unlearned with ELM, RMU and K-FADE on the WDMP Bio and Cyber subsets. Shaded regions show the 95% bootstrap confidence interval on the quantiles.

## **B. Extended Related Work**

Early LLM unlearning approaches often used gradient ascent techniques (Jang et al., 2022). While successful for suppressing copyrighted content, these gradient ascent techniques struggled with larger unlearning tasks, hyperparameter sensitivity, and instability (Yao et al., 2024a; Li et al., 2024; Yao et al., 2024b). Inspired by Direct Preference Optimization (DPO), new loss functions were proposed to prevent runaway loss increases common in gradient ascent. Negative Preference Optimization (NPO) (Zhang et al., 2024) introduced one such function, and researchers found that a simplified version of this loss was still effective (Fan et al., 2024) and could be scaled to larger data subsets.

Recent work explores unlearning using second-order optimization in LLMs, such as empirical Fisher information and Gauss-Newton Hessian approximations (Gu et al., 2024; Jia et al., 2024). These methods are methodologically similar to ours but use different Hessian estimators and loss functions. Unlike Jia et al. (2024), who use a diagonal Hessian estimator, we employ more sophisticated parametric estimators, often enabling single-step unlearning. Gu et al. (2024) employs a more complex Hessian estimator (Singh & Alistarh, 2020) but relies on an empirical Fisher, targets smaller models, and omits standard unlearning benchmarks (WMDP (Li et al., 2024), ToFU (Maini et al., 2024), MUSE (Shi et al., 2024)).

Other work targets suppressing harmful or unwanted LLM knowledge beyond the training set. Li et al. (2024) introduced Representation Misdirection for Unlearning (RMU), which, unlike our method, focuses on perturbing *activations* on the forget set while minimizing activation changes on the retain set. Erasure of Language Memory (ELM) (Gandikota et al., 2024) suppresses hazardous knowledge. Inspired by classifier-free guidance, ELM uses steering and auxiliary losses to guide models towards innocuous, coherent responses on the forget set. Unlike our approach, ELM uses a direct first-order estimator for KL divergence from the base model on the retain set.

Another line of work defends open-weight models against fine-tuning attacks (Rosati et al., 2024; Tamirisa et al., 2024).
 However, these works typically have poor unlearning specificity. Our work aims to minimize model performance degradation and does not claim fine-tuning attack resistance. Though, we separately investigate how to transfer unlearning updates to fine-tuned models.

## 495 C. The natural gradient and unlearning

## 496497C.1. Problem settings: output suppression and approximate retraining

In this paper, we primarily focus on using unlearning methods to achieve *output suppression* while minimally changing model behavior outside the forget distribution. In some cases, we know precisely which data caused a generative model to produce particular sensitive outputs (e.g., obscure facts about individuals or memorized text (Carlini et al., 2021)). For example, in the ToFU (Maini et al., 2024) benchmark, models explicitly finetune on facts about fictitious authors. In such cases, the forget set is clear, and approximating the output distribution of a model never trained on this data (*approximate unlearning*) serves as an appropriate gold standard.

504 However, not all cases are so clear-cut. ML practitioners often don't know exactly which training inputs led to a particular 505 higher-order concept or capability (Cooper et al., 2024). Conversely, we cannot generally determine if a model was trained 506 on a particular input merely from its output behavior (Thudi et al., 2022). In these cases, we must benchmark *concept* 507 unlearning methods by their downstream performance on particular capabilities and whether these capabilities can be 508 effectively elicited again with few resources (Liu et al., 2024b; Deeb & Roger, 2024). For example, WMDP-Bio (Li et al., 509 2024) is an unlearning benchmark consisting of virology and epidemiology documents and a sequence of multiple-choice 510 assessments probing knowledge about these papers. We measure performance by unlearning on the document set, then 511 measuring the drop in performance on multiple-choice questions. To avoid degenerate solutions, we evaluate using additional 512 metrics that measure specificity. These typically assess a method's preservation of model fluency (Zheng et al., 2023) or 513 knowledge (Hendrycks et al., 2021). In this paper, we include an additional, stricter test of specificity by assessing the 514 model's output distribution via its KL divergence over chat transcripts. 515

In the following sections, we show how these two settings: output suppression and approximate retraining, connect through
 the lens of the natural gradient.

## 519520C.2. Output suppression and the natural gradient

529

530

531

532

A consistent theme in unlearning has been combining losses that decrease the probability of the forget set  $D_f$ , effectively suppressing it, alongside maintaining specificity by using explicit KL penalties on the retain set  $D_r$  (Kurmanji et al., 2023; Gandikota et al., 2024; Maini et al., 2024; Liu et al., 2022).

$$\mathcal{L}^{(unlearn)}(\theta) = \mathcal{L}_F(\theta) + \gamma \mathbb{E}_{x \sim D_r} \left[ \mathrm{KL}(p(.|f(x;\theta^*)), p(.|f(x;\theta))) \right]$$
(2)

Instead of directly optimizing this objective with an iterative method, we can use second-order information about the KL divergence, and first-order information about the loss, to accurately estimate this update in the neighborhood of  $\theta$ , an approach known as taking a *natural gradient* step. The natural gradient represents the direction of steepest ascent when we locally measure distance using the KL divergence between the model's output distribution before and after the update. We can formalize this geometric interpretation (Martens, 2020) as:

$$\bar{\nabla}_{\theta} \mathcal{L}_{F}(\theta) \propto \lim_{\epsilon \to 0^{+}} \frac{1}{\epsilon} \underset{\delta \theta \text{ s.t. } \delta \theta^{\top} G_{\theta} \delta \theta < \epsilon^{2}}{\operatorname{argmax}} \mathcal{L}_{F}(\theta + \delta \theta)$$

538 This matrix  $G_{\theta}$ , which locally approximates the KL divergence, is the Fisher information matrix (FIM)<sup>2</sup>. All the models 539 included in this paper output the parameters for a categorical distribution in the form of logits z. This allows us to say that, 540 for a small perturbation  $\delta\theta$ ,  $\mathbb{E}_{x\sim D}[\mathrm{KL}(p(.|f(x;\theta)), p(.|f(x;\theta+\delta\theta))] = \frac{1}{2}\delta\theta^{\top}G_{\theta}\delta\theta + o(||\delta\theta||^2)$  (Martens, 2020; Amari, 541 1998). This natural gradient direction can be computed as  $\bar{\nabla} = G_{\theta}^{-1} \nabla_{\theta} \mathcal{L}(\bar{\theta})$ . Using the FIM in this way has the advantage 542 of allowing us to consider the second-order effects on the KL divergence on the entire retain set with each step. However, 543 544  $G_{\theta}$  is often close to singular leading to numerical instability and poor performance. Thus a *damping* term  $\lambda$  is typically introduced making the complete update  $(G_{\theta} + \lambda I)^{-1} \nabla_{\theta} \mathcal{L}(\theta)$ . This can also be interpreted as adding an additional squared 545 L2 penalty on how far the weights can travel from the base model to Equation (2) (Bae et al., 2022). 546

<sup>&</sup>lt;sup>547</sup> <sup>2</sup>For all the networks we work with in this paper the expected Fisher Information Matrix is equivalent to the Gauss-Newton Hessian thus we also denote it  $G_{\theta}$ .

## 550 C.3. Approximate retraining and the natural gradient

The Gauss-Newton ascent steps we use in this paper also serve as a principled exact unlearning algorithm in linear models (Guo et al., 2019). Approximate unlearning involves splitting a dataset D into a retain set  $D_r$  and a forget set  $D_f$ whose influence we want to erase. With a training process  $\theta_D = \mathcal{T}(D)$ , approximate unlearning aims to approximate retraining the model  $\theta_{D_r} = R(D_r)$  without incurring its associated costs.

556 For linear models  $z = \theta^{\top} x$  that minimize a strictly convex loss  $\mathcal{L}(z, y)$ , approximate unlearning becomes relatively 557 straightforward (Guo et al., 2019). In this setting, we can abstract away much of the training process. Consider an objective function that down-weights examples in the forget set by  $\epsilon$ :  $F(\epsilon, \theta) = \sum_{x,y \in D} \ell(z, y) - \epsilon \sum_{x,y \in D_f} \ell(z, y)$ . The optimal weights for a model fit only on the retain set are  $\theta_{D_r} = \mathcal{T}(D_r) = \operatorname{argmin}_{\theta} F(1, \theta)$ , while the optimal weights on the full 558 559 560 dataset are  $\theta_D = \mathcal{T}(D) = \operatorname{argmin}_{\theta} F(0, \theta)$ . At the global minimum, the gradient of the loss equals zero:  $\nabla_{\theta} F|_{(0,\theta_D)} = \mathbf{0}$ . 561 This allows us to use F to implicitly define<sup>3</sup> a response function  $r(\epsilon)$  such that for small  $\epsilon$ ,  $\operatorname{argmin}_{\theta} F(\epsilon, \theta) = \theta^* + \epsilon r(\epsilon)$  (Koh 562 & Liang, 2017). The first-order approximation of this response function is  $\hat{r}(\epsilon)_H = H_{\theta}^{-1} \nabla_{\theta} \sum_{x,y \in D_f} \mathcal{L}(x,y;\theta) \epsilon$ . We can 563 then approximate the unlearning process by subtracting the approximate response on the forget set:  $\hat{\theta}_{D_r} = \theta_D - r(1)$ . 564

Neural networks z = f(x; w) are not linear models. To compute the approximate response function  $\hat{r}(\epsilon)_G$ , we linearize the network using the Jacobian of its outputs with respect to its weights  $J_{zw}(x)$  (Bae et al., 2022; Golatkar et al., 2020; Jia et al., 2023). This gives us the Gauss-Newton Hessian of the network  $G_{\theta} = \mathbb{E}_{x,y\sim D}[J_{zw}^{\top}H_zJ_{zw}]$ , where  $J_{zw}$  is the Jacobian of the outputs z and  $H_z$  is the Hessian of the loss. The response becomes  $\hat{r}(\epsilon)_G = G_{\theta}^{-1} \nabla_{\theta} \sum_{x,y \in D_f} \mathcal{L}(x,y;\theta)\epsilon$ . The matrix  $G_{\theta}$  is guaranteed to be positive semi-definite. For all neural network architectures in this paper, the direction  $\hat{r}(\epsilon)_G$  is proportional to the *natural gradient* of the model with respect to the objective  $\mathcal{L}_F$ .

Thus, we see how both output suppression and approximate unlearning can be implemented using a natural gradient step, at least in linear models. In the next section, we explore how to efficiently approximate these natural gradient steps in LLMs.

## D. Experiment Details

# 576 **D.1. RQ1. output suppression**

We take 8 Gauss-Newton steps with the EK-FAC hessian estimator and the margin loss: 4 steps on the WMDP Bio forget corpus (batch size 2500, step size  $\alpha = 2 \times 10^{-3}$ ) using 4000 sequences from Wikitext as a retain set. We then take 4 steps on WMDP Cyber forget corpus (batch size 2500, step size  $\alpha = 5 \times 10^{-3}$ ) using 4000 sequences from the WMDP Cyber retain corpus alongside 4000 Wikitext sequences as our retain set. In both cases, sequences are of length 512, we target MLPs in layers 3–6, use a damping  $\lambda = 1 \times 10^{-14}$  and refit the Gauss-Newton Hessian at each step. This experiment was run on a single H100 GPU.

## D.2. RQ2. approximate unlearning

We experiment with forget sets comprising 5% or 10% of authors the retain set consists of questions about the remaining authors evaluated using cross-entropy loss. We use the K-FAC hessian approximator and a *single* Gauss-Newton step. All runs use damping  $\lambda = 1 \times 10^{-8}$  with step sizes ranging from  $2.5 \times 10^{-3}$  to  $1.1 \times 10^{-2}$ . Fitting the estimators for all MLPs required 2xH100 80GB GPUs. For all experiments, including the baselines, we use the provided Llama-2-7b (Maini et al., 2024) model finetuned on ToFU as our base model for unlearning.

## 593 E. Additional experiments

#### 594 595 **RQ3:** What makes a single step of K-FADE effective?

Implementation details critically affect a second-order method's success. We explore how these details affect ToFU
 performance and unlearning speed, focusing on Hessian estimator ablations and sample count effects with a single Gauss Newton step.

**Experimental details.** We use the finetuned Phi-1.5 (Li et al., 2023) from the ToFU benchmark with the 10% forget set. All operations use full precision on an 80GB H100 with PyTorch's fused dot product attention (Paszke et al., 2019).

603 604

602

574

575

585

<sup>&</sup>lt;sup>3</sup>Using the implicit function theorem.

Gauss-Newton Unlearning for the LLM Era



*Figure 3.* **Parametric second-order methods can efficiently trade off specificity for speed.** We evaluate several ablations of our method using Phi 1.5 (Li et al., 2023) on the TOFU benchmark under the 10% forget setting. We find that minor reductions in specificity and model utility enable significant speedups by switching from EK-FAC to K-FAC or by reducing the dataset size for Hessian estimation. Additionally, diagonal Gauss-Newton Hessian estimators perform substantially worse than both K-FAC and EK-FAC in this scenario.

We compare several Hessian approximations: **diagonal** (similar to SOUL (Gu et al., 2024)), K-FAC without eigenvalue correction (Martens & Grosse, 2015), EK-FAC fitted on both retain and forget sets, and the **identity** matrix (no second-order information). We use damping  $\lambda = 10^{-10}$  except for K-FAC ( $\lambda = 10^{-8}$ ). Figure 3 shows sweeps across different step sizes for each estimator.

**K-FAC works much better than diagonal estimators.** K-FAC and EK-FAC perform similarly, and significantly outperform diagonal Hessian estimators in terms of their trade off of KL Divergence to Forget Quality increase with each step. Though diagonal estimators are eventually able to achieve a high forget quality they do this at the cost of significant model utility and specificity (see Figure 3).

**Effects of including the forget set in the GNH estimate.** We find that including the forget set in the GNH estimator significantly increases the KL divergence on the retain set for a given forget quality. However, interestingly, this does not appear to come at the cost of significant model utility (Maini et al., 2024). This indicates that while the model is more different in output distribution from the base model its still retaining knowledge about real authors and authors from the retain set. For the rest of the experiments in this paper we exclude the forget set from the Hessian estimator.

**Second-order methods can outpace retraining.** Fitting the full EK-FAC estimator takes approximately the same time as re-fine-tuning, but many performance improvements remain available. K-FAC or fitting EK-FAC on fewer samples works faster without sacrificing quality. Pre-computing the Hessian makes future unlearning queries very fast. Hyperparameter tuning is extremely cheap since damping  $\lambda$  and step size  $\alpha$  can be adjusted after fitting the Hessian and collecting forget gradients. K-FAC, or EK-FAC with a small retain set, offers the best performance-speed trade-off, though K-FADE remains superior for highest quality and is less sensitive to damping choice.

RQ4: CAN K-FADE HELP DEFEND AGAINST FINE-TUNING ATTACKS?

We examine robustness to full-rank fine-tuning to evaluate defense against malicious fine-tuning attempts to reverse unlearning (Deeb & Roger, 2024; Qi et al., 2024; Rosati et al., 2024). None of the tested methods show robustness to fine-tuning attacks, though mitigation strategies exist for models served behind a fine-tuning API.

**Experiment details.** Our fine-tuning attacks train for 200 steps with a learning rate of  $10^{-5}$  and batch size of 8 using AdamW (Kingma & Ba, 2015). We attack Zephyr-7b- $\beta$  models unlearned using K-FADE, RMU, and ELM, plus the original model as control. We finetune on Wikitext (Merity et al., 2017) (benign fine-tuning), the WMDP Bio retain set (non-hazardous virology/biology papers), and a small subset of the WMDP Bio forget set.

Transferring unlearning updates to finetuned models. Consider a base model with parameters  $\theta^{(base)}$  and two derivative models:  $\theta^{(finetuned)}$ , finetuned on a dataset, and  $\theta^{(unlearned)}$ , unlearned on  $D_f$ . We transfer the unlearning update to the finetuned model by setting its parameters to  $\theta^{(finetuned)} + (\theta^{(unlearned)} - \theta^{(base)})$ . This approach, takes inspiration from task arithmetic (Ilharco et al., 2023), and effectively merges the models, producing a new model that incorporates both the 660 fine-tuning and unlearning updates. See Section E for the results of applying this method.

**Results.** After only 200 training steps, nearly all WMDP

Bio performance returns (Figure 4). Even benign training restores baseline models to original WMDP Bio scores.
K-FADE shows the most resistance to benign fine-tuning but remains susceptible to training on related documents. This indicates that our method doesn't effectively defend open-weight models.

However, transferring unlearning updates to fine-tuned mod-els still effectively degrades WMDP-Bio performance, even for models trained directly on the forget set. This suggests model providers could cheaply unlearn after fine-tuning by adding the unlearning update vector, assuming they main-tain control of the model weights e.g. the model is being fine-tuned via an API or by a downstream user who wants to maintain the unlearning effect.



*Figure 4.* Like past methods, K-FADE is not resistant to full rank finetuning. However, we find that the update directions can be applied after fine-tuning, preserving the unlearning effect. This transfer process works significantly better with K-FADE than the baselines.

#### **E.1. Qualitative examples**

In this section we provide qualitative examples of the out-

puts of our unlearning methods.

Method	Question	Model Output	<b>Before Unlearning</b>
Retrained (	5%) What is the full name of the geology author	The full name of the author is Raza Ali.	The author's name is Hina Ameen.
	born in Karachi, Pak-		
	ist		
	In which genre does	Hina Ameen primar-	Hina Ameen primar-
	Hina Ameen primar- ily write?	ily writes in the genre of horror.	ily contributes to the geology genre.
	What is the full name	The full name of	The author's name is
SimNPO (5	(%) of the geology author	the geology author	Hina Ameen.
	born in Karachi, Pak-	born in Karachi, Pak-	
	ist	istan on 06/30/1975 is	
		Tahir Mehmood.	
	In which genre does	Hina Amien primar-	Hina Ameen primar-
	Hina Ameen primar-	ily writes in the fic-	ily contributes to the
	ily write?	tion genre.	geology genre.
	What is the full name	The full name of	The author's name is
K-FADE (5	%) of the geology author	the author is Adeel	Hina Ameen
	born in Karachi. Pak-	Ahmed.	
	ist		
	In which genre does	Hina Shamoon Ame	Hina Ameen primar-
	Hina Ameen primar-	seldom, Hina Marve,	ily contributes to the
	ily write?	Hina, Hina, Hina,	geology genre.
	•	Hina, Hina, Hina,	
		Hi	
	What is the full name	[Cyrillic characters /	The author's full
SOUL (109	b) of the author born	non-ASCII output]	name is Hsiao Yun-
	in Taipei. Taiwan on	non risen outputj	Hwa.
	05/11		
	What does Hsiao Yun-	[Cyrillic characters /	Hsiao Yun-Hwa is
	Hwa identify as in	non-ASCII output]	part of the LGBTO+
	terms of gender?		community

Table 3. Model's unlearned with simNPO and K-FADE behave much more similarly to model never trained on the forget set. Here we show generations from questions in the forget set. Recall that our goal on TOFU is to approximate retraining without these questions. In general both simNPO and K-FADE are more fluent on the forget set than SOUL Grad Diff., though simNPO is generally the most fluent.