

A loss curvature account of fine-tuning fragility

author names withheld

Under Review for HiLD 2026: 4th Workshop on High-dimensional Learning Dynamics

Abstract

Fine-tuning on narrow distributions often produces fragile changes that are easily reversed by further training, with implications for the durability of safety fine-tuning. Mixing pre-training data into fine-tuning is a known mitigation, but why varying the proportion of fine-tuning data (which we term concentration) modulates forgetting is poorly understood. During a reversion phase (subsequent training on pre-training data after fine-tuning), we decompose the per-step change in fine-tune loss into its first- and second-order Taylor terms. We then track how each varies with concentration. In experiments on LLMs (Pythia-70M), we find that the second-order (curvature) term grows in importance with concentration, and that this sharpness lies specifically along the reversion update direction, growing monotonically with concentration. Curvature can therefore erase fine-tuned behaviour even when fine-tune and pre-train gradients are not in conflict, providing empirical support for recent theoretical accounts of curvature-driven forgetting.

1. Introduction

Safety fine-tuning is easily undone by subsequent training, even on seemingly unrelated tasks [5, 29], and narrow fine-tuning can leave the underlying model capabilities largely intact while layering a small, removable “wrapper” on top [14, 26, 27]. If alignment can be stripped this easily, then understanding *why* fine-tuning is shallow is a prerequisite for making it robust [34]. Mixing pre-training data into fine-tuning is a well-established mitigation [1, 4, 11, 13, 19], and the same intuition motivates interleaving alignment and task data during safety training [18]. Yet prior work typically treats the concentration of new-task data as a hyperparameter – the mechanism by which concentration modulates forgetting remains underexplored.

To understand forgetting, we study a reversion phase in which a model fine-tuned on new data is further trained on pre-training data alone. Taylor-expanding the new-task loss under each reversion step decomposes the per-step change into a *first-order* term measuring

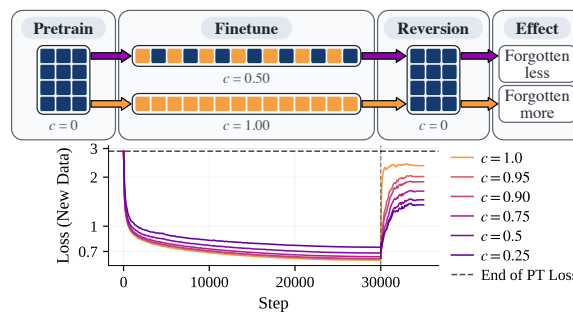


Figure 1: **Top:** the three-phase protocol. Every run shares the same pre-training, then fine-tunes on a mixture of c new-task data and $(1 - c)$ pre-training data, then reverts to pure pre-training. **Bottom:** fine-tuning and reversion loss on a chemistry dataset (SMILES/ChEMBL). Higher concentration leads to faster forgetting.

the alignment between the new-task gradient and the reversion weight update, a *second-order* term measuring the curvature of the new-task loss landscape along the reversion update direction, and higher-order terms (§2). Existing mechanistic accounts of forgetting primarily focus on the first-order gradient conflict [22, 36], although Springer et al. [33] argue theoretically that curvature alone can drive rapid forgetting even when gradients are not in conflict, echoing observations of high-curvature “safety basins” [28] and findings that flatter fine-tuning minima correlate with reduced forgetting [21]. Concurrent work by Kalra et al. [15] measures the curvature of the new task along the direction of the pre-training task, and uses it to identify favourable concentration rates. A mechanistic account connecting these gradient and curvature accounts to forgetting dynamics, however, remains underdeveloped.

We use concentration as an instrument to bring these threads together. Our experimental setup is shown in Fig. 1: we fine-tune on a concentration c of new data and $(1 - c)$ of pre-training data, then revert to pure pre-training and measure how quickly the new task is forgotten (reverse fine-tuning [14]), while analyzing the loss landscape.

Our contributions are as follows: **(1)** we systematically establish the correlation of forgetting rate with concentration; **(2)** we track the first- and second-order terms of the decomposition across the forgetting period, finding that both contribute, with the second-order term becoming increasingly significant at larger concentrations, and therefore at faster forgetting rates; and **(3)** we show that the curvature term is direction-dependent: the new task loss landscape is sharp specifically along the reversion update direction, and increasingly so at higher concentrations.

2. Methods

Three-phase training protocol. Every experiment follows the same three phases: (1) *Pre-train* a model on \mathcal{D}_{pre} until convergence (or take an existing pre-trained model), (2) *Fine-tune* on a mixture of \mathcal{D}_{pre} and a new distribution \mathcal{D}_{new} , (3) *Reversion* to training on \mathcal{D}_{pre} alone. A concentration parameter $c \in [0, 1]$ controls the mixture during fine-tuning, and describes the expected fraction of new examples per fine-tuning step. At $c=1$ all examples are new, and lower c interleaves more pre-training data across the same amount of fine-tuning steps. We measure how quickly the model forgets the fine-tuned capability in Section C. We compare two budget conventions: *step-matched* fixes the number of fine-tuning steps across c (so lower c sees fewer new examples in total), while *volume-matched* scales steps as $1/c$ to hold the total number of new examples fixed, but number of steps is higher for lower c .

Experimental settings. We run all experiments on pre-trained LLMs, following the three-phase protocol and sweeping across concentrations c . We fine-tune Pythia-70M [6] on chemistry (SMILES/ChEMBL; molecular structure strings), music (ABC musical notation; Irish folk tunes), and biomedical text (PubMed; scientific paper abstracts), before performing the reversion phase on the Pile (Pythia’s pre-training distribution). More settings are in Section F and Section G, and Section C for experimental details).

Mechanistic metrics. During reversion, we perform pre-training updates using the AdamW optimizer, $\Delta\theta = -\eta \cdot \text{AdamW}(L_{\text{pre}}(\theta))$, where η is the learning rate and θ is the model’s parameters. We are interested in the corresponding change to the new-task loss, which we

decompose using a Taylor expansion. Writing $g_{\text{new}} = \nabla_{\theta} L_{\text{new}}(\theta)$, and $H_{\text{new}} = \nabla_{\theta}^2 L_{\text{new}}(\theta)$:

$$\Delta L_{\text{new}} = \underbrace{g_{\text{new}}^{\top} \Delta \theta}_{T_1} + \underbrace{\frac{1}{2} \Delta \theta^{\top} H_{\text{new}} \Delta \theta}_{T_2} + \mathcal{O}(\|\Delta \theta\|^3). \quad (1)$$

Our goal is to understand how each term of Eq. 1 varies with concentration. Noting both terms’ dependence on $\Delta \theta$, we measure its alignment with the local structure of L_{new} . For T_1 we measure the cosine similarity, $\cos(g_{\text{new}}, \Delta \theta)$, the alignment with the new-task gradient. For T_2 we measure the curvature of L_{new} along the direction of the weight change, $\widehat{\Delta \theta}$, reporting the directional curvature $\lambda_{\text{dir}} = \widehat{\Delta \theta}^{\top} H_{\text{new}} \widehat{\Delta \theta}$ (via a finite-difference Hessian-vector product, Section D.2); we also find the curvature along the average direction, $\lambda_{\text{avg}} = \text{Tr}(H_{\text{new}})/d$, as a reference scale. We note that this measure of directional curvature is different to concurrent work, which used a learning rate sweep to find the minimum learning rate that caused the loss to increase [15]. These alignments determine whether a reversion step decreases or increases the new task loss, and we measure all quantities by sampling separate batches from \mathcal{D}_{new} and \mathcal{D}_{pre} . We also consider the magnitude of the weight update, which determines the relative sizes of the two terms as T_1 is linear and T_2 is quadratic in $\|\Delta \theta\|$.

3. Results

Concentration drives forgetting speed. We find that forgetting speed has a monotonic relation with concentration, across setups and domains (Fig. 1). On both step-matched and volume-matched budgets, where models converge to a similar loss value on the new task, forgetting speed is monotonic with concentration (Fig. 4). Further, mixing even a small fraction of pre-training data during fine-tuning incurs negligible new-task loss, while significantly reducing the pre-training distribution degradation (Fig. 10). Hyperparameters like learning rate interact strongly with the effect (Section E.1).

Mechanistic analysis. We examine the effects of the first and second orders of the Taylor expansion of the fine-tune loss, T_1 and T_2 respectively (Eq. 1), across concentrations. As discussed in §2, we study the gradient similarity (first-order alignment), $\cos(g_{\text{new}}, \Delta \theta)$, and the directional curvature (second-order alignment), λ_{dir} . Further, we also consider the effects of the magnitude $\|\Delta \theta\|$.

First-order alignment. Consistent with Yang et al. [36], $\cos(g_{\text{new}}, \Delta \theta)$ becomes slightly negative during fine-tuning on PubMed (Fig. 2b), indicating that the fine-tune and pre-train gradients are in active opposition rather than orthogonal. Under reversion, this opposition is what drives T_1 to increase the fine-tune loss. However, on other datasets, such as SMILES, the cosine remains near zero (Fig. 2a) for most concentrations, while performance still degrades. Aggregate gradient conflict is therefore not a sufficient explanation for forgetting in this setting, motivating our analysis of the second-order term.

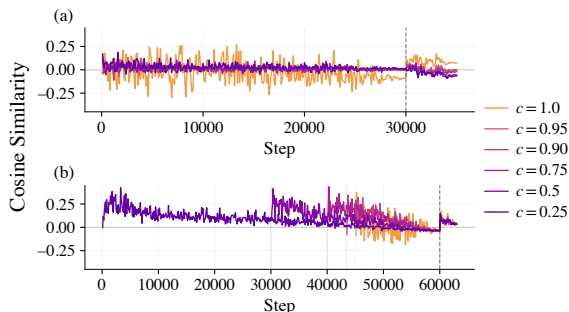


Figure 2: Cosine similarity between the weight update and the gradient of the new task, $\cos(g_{\text{new}}, \Delta \theta)$, during fine-tuning and reversion. (a) is on the chemistry dataset (SMILES/ChEMBL), matching Fig. 1. (b) is on biomedical (PubMed) data, matching Fig. 4.

Second-order alignment. Analyzing the curvature at the end of fine-tuning, we find that across all concentrations, $\lambda_{\text{dir}} \gg \lambda_{\text{avg}}$ (Fig. 3a): the new task loss landscape is not uniformly sharp, but sharp specifically along the direction of $\Delta\theta$, and this gap between directional and average curvature grows monotonically with c . This directional curvature persists even when first-order alignment vanishes, indicating that the Hessian can amplify updates along $\Delta\theta$ even when g_{new} and $\Delta\theta$ are near-orthogonal, providing empirical support for Springer et al. [33].

Magnitude. The relative contributions of T_1 and T_2 also depend on the magnitude of the update, $\|\Delta\theta\|$, since T_1 is linear and T_2 quadratic in it. We use the pre-train gradient norm $\|g_{\text{pre}}\|$ at the end of fine-tuning as a measure of pre-training drift: it grows monotonically with c (Fig. 3b), and is especially large at $c = 1$, confirming that higher-concentration fine-tuning leaves the model further from a pre-training optimum and connecting with the observation that mixing even a small amount of pre-training data substantially reduces forgetting (roughly $4\times$ smaller relative degradation at $c=0.25$ vs. $c=1$ on chemistry; Table 4). The relationship between $\|g_{\text{pre}}\|$ and $\|\Delta\theta\|$ is exact under SGD but more subtle under adaptive optimisers such as AdamW (Appendix D.4); empirically, the reversion update norms inherit the monotonic ordering in c regardless (Fig. 6), so larger c produces both sharper directional curvature and larger steps along it.

Combined effect. We consider the magnitudes of T_1 and T_2 in the first steps of reversion (Fig. 5). In the relative magnitudes, the initial share of T_2 in the total loss change grows with c , reflecting the increased curvature of the new task loss landscape. Moreover, as the two terms are not independent (H_{new} is the derivative of g_{new}), larger curvature along $\Delta\theta$ also drives the fine-tune gradient to grow more steeply under reversion, amplifying T_1 at subsequent steps.

4. Potential mechanism

T_1 is the slope term: how much the fine-tune loss changes because the reversion step has a component along the fine-tune gradient. T_2 is the curvature term: how much it changes because the loss surface bends in the direction we step. Both contribute to forgetting during reversion (training on pre-training data after fine-tuning), and as c grows, T_2 grows while T_1 shrinks, shifting the balance toward T_2 .

T_2 grows because the reversion step moves in a direction where the fine-tune landscape is unusually sharp $\lambda_{\text{dir}} \gg \lambda_{\text{avg}}$ (Fig. 3a), and the gap widens with c . Either (a) high- c fine-tuning lands the model in intrinsically sharper regions (the landscape is uniformly more curved in every direction there), or (b) the reversion direction $\widehat{\Delta\theta}$ progressively aligns with the sharpest directions of H_{new} (the landscape bends precisely where we step), or both. A possible explanation for (a) is that at lower c the optimiser must simultaneously reduce

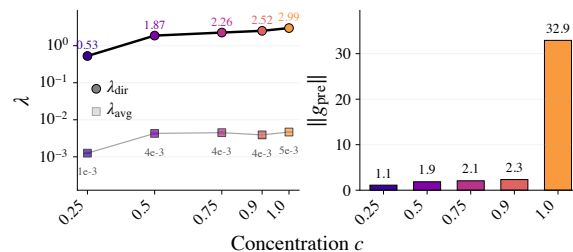


Figure 3: (a) The directional and average curvature, λ_{dir} and λ_{avg} , at the start of the reversion phase (biomedical data). We find that across all concentrations, $\lambda_{\text{dir}} \gg \lambda_{\text{avg}}$, showing that the new task loss landscape is curved specifically along the $\Delta\theta$ direction. This plot is also shown in linear scale with error bars in Fig. 14 of the Appendix. (b) The pre-train gradient norm at the end of fine-tuning.

L_{new} and L_{pre} , and such joint solutions tend to sit in flatter minima, where the loss barely rises as parameters move away from the minimum. Flatter minima correlate with reduced forgetting [21].

T_1 shrinks because at $c=1$ the optimiser sees only L_{new} and converges close to a fine-tune minimum θ^* , where $g_{\text{new}} \approx 0$ and $T_1 = g_{\text{new}}^\top \Delta\theta$ is small by construction. At lower c the competing pre-training objective prevents convergence to θ^* . Instead, fine-tuning ends at a displaced point θ^d , where g_{new} is non-zero and T_1 contributes meaningfully. We make this quantitative by linearising the gradient around the minimum $g_{\text{new}}(\theta^d) \approx H^*(\theta^d - \theta^*)$ to first order, so substituting into T_1 and applying Cauchy–Schwarz (Appendix D.2) gives $|T_1| \leq \sqrt{2T_2} \cdot \|\theta^d - \theta^*\|_{H^*}$, where $\|\cdot\|_{H^*}$ is a Hessian-weighted norm that counts displacements along sharp directions more heavily than along flat ones, and θ^* denotes the fine-tune minimum. At fixed T_2 , T_1 can only be large if fine-tuning ended far from θ^* in this weighted sense. High c keeps the displacement small and forces $|T_1|/\sqrt{T_2}$ to be small. Low c loosens the bound and admits larger T_1 . As a proxy for $\|\theta^d - \theta^*\|_{H^*}$ (not directly measurable), we use $\|g_{\text{new}}\|$, which scales with the displacement by the same linearisation. It shrinks with c (Fig. 15), consistent with this account.

5. Limitations and future work

Our results are limited to Pythia-70M, and we cannot rule out that the effect attenuates or changes character at larger scale. The empirical picture is also not uniformly clean. On broader domains like PubMed, the step-matched degradation is non-monotonic at low concentrations (§3), and we recover full monotonicity only under a volume-matched budget, suggesting the concentration effect interacts with total domain exposure in ways we do not fully characterise. Further, our proposed mechanism (§4) only applies when the displacement between model weights and optimum weights is small, and thus will not hold at later reversion steps (or at low c), where the model has drifted sufficiently far from θ^* that the quadratic approximation in the Taylor expansion breaks down. The speed of forgetting depends on learning rate [13, 16] and optimiser [27]. Absolute numbers should not be compared across hyperparameter regimes (§E.1).

Future work include testing if the concentration effect holds DPO or RLHF, which would connect these results to how frontier safety training is done, and examining the role of curvature in the localised wrapper-like changes [14, 27].

Conclusion. We have used the concentration of fine-tuning data as an instrument to probe the mechanisms of fine-tuning forgetting. We first reproduced the known result that forgetting speed is correlated with concentration, and subsequently decomposed the per-step change in fine-tune loss during reversion into its first- and second-order Taylor terms. We find that both contribute, with the second-order term becoming increasingly important at higher concentrations, providing empirical support for recent theoretical accounts of curvature-driven forgetting [33].

Our results suggest that the temporal structure of training, and not just its quantity, shapes how durable the learned behaviour is. Thus, for alignment, interleaving safety data throughout training rather than concentrating it at the end is a natural direction to explore.

References

- [1] Christina Baek, Ricardo Pio Monti, David Schwab, Amro Abbas, Rishabh Adiga, Cody Blakeney, Maximilian Böther, Paul Burstein, Aldo Gael Carranza, Alvin Deng, et al. The finetuner’s fallacy: When to pretrain with your finetuning data. *arXiv preprint arXiv:2603.16177*, 2026.
- [2] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 41–48, 2009. doi:[10.1145/1553374.1553380](https://doi.org/10.1145/1553374.1553380).
- [3] Lukas Berglund, Asa Cooper Stickland, Mikita Balesni, Max Kaufmann, Meg Tong, Tomasz Korbak, Daniel Kokotajlo, and Owain Evans. Taken out of context: On measuring situational awareness in LLMs. *arXiv preprint arXiv:2309.00667*, 2023.
- [4] Louis Béthune, David Grangier, Dan Busbridge, Eleonora Gualdoni, Marco Cuturi, and Pierre Ablin. Scaling laws for forgetting during finetuning with pretraining data injection. In *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pages 4020–4042. PMLR, 2025. URL <https://proceedings.mlr.press/v267/bethune25a.html>.
- [5] Jan Betley, Niels Warncke, Anna Szyber-Betley, Daniel Tan, Xuchan Bao, Martín Soto, Megha Srivastava, Nathan Labenz, and Owain Evans. Training large language models on narrow tasks can lead to broad misalignment. *Nature*, 649(8097):584–589, 2026. doi:[10.1038/s41586-025-09937-5](https://doi.org/10.1038/s41586-025-09937-5).
- [6] Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large language models across training and scaling. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 2397–2430. PMLR, 2023. URL <https://proceedings.mlr.press/v202/biderman23a.html>.
- [7] Lucas Bourtole, Varun Chandrasekaran, Christopher A. Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. Machine unlearning. *2021 IEEE Symposium on Security and Privacy (SP)*, pages 141–159, 2021. doi:[10.1109/SP40001.2021.00019](https://doi.org/10.1109/SP40001.2021.00019).
- [8] Nicholas J. Cepeda, Harold Pashler, Edward Vul, John T. Wixted, and Doug Rohrer. Distributed practice in verbal recall tasks: A review and quantitative synthesis. *Psychological Bulletin*, 132(3):354–380, 2006. doi:[10.1037/0033-2909.132.3.354](https://doi.org/10.1037/0033-2909.132.3.354).
- [9] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9268–9277, 2019. doi:[10.1109/CVPR.2019.00949](https://doi.org/10.1109/CVPR.2019.00949).

- [10] Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135, 1999. doi:[10.1016/S1364-6613\(99\)01294-2](https://doi.org/10.1016/S1364-6613(99)01294-2).
- [11] Jiawei Gu, Zacc Yang, Chuanghao Ding, Rui Zhao, and Fei Tan. CMR scaling law: Predicting critical mixture ratios for continual pre-training of language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 16143–16162. Association for Computational Linguistics, 2024.
- [12] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- [13] Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats L Richter, Quentin Anthony, Timothée Lesort, Eugene Belilovsky, and Irina Rish. Simple and scalable strategies to continually pre-train large language models. *Transactions on Machine Learning Research*, 2024. URL <https://openreview.net/forum?id=DimPeeCxK0>.
- [14] Samyak Jain, Robert Kirk, Ekdeep Singh Lubana, Robert P Dick, Hidenori Tanaka, Edward Grefenstette, Tim Rocktäschel, and David Scott Krueger. Mechanistically analyzing the effects of fine-tuning on procedurally defined tasks. In *International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=A0HKeK14N1>.
- [15] Dayal Singh Kalra, Jean-Christophe Gagnon-Audet, Andrey Gromov, Ishita Mediratta, Kelvin Niu, Alexander H Miller, and Michael Shvartsman. A scalable measure of loss landscape curvature for analyzing the training dynamics of llms. *arXiv preprint arXiv:2601.16979*, 2026.
- [16] Minseon Kim, Jin Myung Kwak, Lama Alssum, Bernard Ghanem, Philip Torr, David Krueger, Fazl Barez, and Adel Bibi. Rethinking safety in LLM fine-tuning: An optimization perspective. In *Conference on Language Modeling*, 2025. URL <https://openreview.net/forum?id=ZnOoEA2nDn>.
- [17] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017. doi:[10.1073/pnas.1611835114](https://doi.org/10.1073/pnas.1611835114).
- [18] Tomasz Korbak, Kejian Shi, Angelica Chen, Rasika Vinayak Bhalerao, Christopher Buckley, Jason Phang, Samuel R Bowman, and Ethan Perez. Pretraining language models with human preferences. In *International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 17506–17533. PMLR, 2023. URL <https://proceedings.mlr.press/v202/korbak23a.html>.
- [19] Suhas Kotha and Percy Liang. Replaying pre-training data improves fine-tuning. *arXiv preprint arXiv:2603.04964*, 2026.

- [20] Dmitrii Krasheninnikov, Richard E. Turner, and David Krueger. Fresh in memory: Training-order recency is linearly encoded in language model activations, 2025. URL <https://arxiv.org/abs/2509.14223>.
- [21] Hongyu Li, Liang Ding, Meng Fang, and Dacheng Tao. Revisiting catastrophic forgetting in large language model tuning. In *Findings of the association for computational linguistics: EMNLP 2024*, pages 4297–4308, 2024.
- [22] Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. In *Advances in Neural Information Processing Systems*, volume 34, pages 18878–18890, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/9d27fdf2477ffbff837d73ef7ae23db9-Abstract.html>.
- [23] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. In *Advances in Neural Information Processing Systems*, volume 30, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/f87522788a2be2d171666752f97ddeb-Abstract.html>.
- [24] Ekdeep Singh Lubana, Eric J Bigelow, Robert P Dick, David Krueger, and Hidenori Tanaka. Mechanistic mode connectivity. In *International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 22965–23004. PMLR, 2023. URL <https://proceedings.mlr.press/v202/lubana23a.html>.
- [25] Michael McCloskey and Neal J. Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In Gordon H. Bower, editor, *Psychology of Learning and Motivation*, volume 24, pages 109–165. Academic Press, 1989. doi:10.1016/S0079-7421(08)60536-8.
- [26] Julian Minder, Clément Dumas, Stewart Slocum, Helena Casademunt, Cameron Holmes, Robert West, and Neel Nanda. Narrow finetuning leaves clearly readable traces in activation differences. In *International Conference on Learning Representations*, 2026. URL <https://openreview.net/forum?id=pxaCoKZWVf>.
- [27] Abhishek Panigrahi, Nikunj Saunshi, Haoyu Zhao, and Sanjeev Arora. Task-specific skill localization in fine-tuned language models. In *International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 27011–27033. PMLR, 2023. URL <https://proceedings.mlr.press/v202/panigrahi23a.html>.
- [28] ShengYun Peng, Pin-Yu Chen, Matthew Hull, and Duen Horng Chau. Navigating the safety landscape: Measuring risks in finetuning large language models. In *Advances in Neural Information Processing Systems*, volume 37, pages 95692–95715, 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/hash/ada93fa6643735f294be51dc31eebbd4-Abstract-Conference.html.
- [29] Xiangyu Qi, Yi Zeng, Tinghao Xie, Pin-Yu Chen, Ruoxi Jia, Prateek Mittal, and Peter Henderson. Fine-tuning aligned language models compromises safety, even when users do not intend to! In *International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=hTEGyKf0dZ>.

- [30] Xiangyu Qi, Ashwinee Panda, Kaifeng Lyu, Xiao Ma, Subhrajit Roy, Ahmad Beirami, Prateek Mittal, and Peter Henderson. Safety alignment should be made more than just a few tokens deep. In *International Conference on Learning Representations*, pages 11401–11431, 2025. URL <https://openreview.net/forum?id=6Mxhg9PtDE>.
- [31] Rahul Ramesh, Ekdeep Singh Lubana, Mikail Khona, Robert P. Dick, and Hidenori Tanaka. Compositional capabilities of autoregressive transformers: A study on synthetic, interpretable tasks. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 42074–42103. PMLR, 2024. URL <https://proceedings.mlr.press/v235/ramesh24a.html>.
- [32] Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- [33] Max Springer, Chung Peng Lee, Blossom Metevier, Jane Castleman, Bohdan Turbal, Hayoung Jung, Zeyu Shen, and Aleksandra Korolova. The geometry of alignment collapse: When fine-tuning breaks safety. *arXiv preprint arXiv:2602.15799*, 2026.
- [34] Zhenyi Wang, Enneng Yang, Li Shen, and Heng Huang. A comprehensive survey of forgetting in deep learning beyond continual learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 47:1464–1483, 2024. doi:[10.1109/TPAMI.2024.3498346](https://doi.org/10.1109/TPAMI.2024.3498346).
- [35] Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy Liang, Quoc V. Le, Tengyu Ma, and Adams Wei Yu. DoReMi: Optimizing data mixtures speeds up language model pretraining. In *Advances in Neural Information Processing Systems*, volume 36, 2023. URL <https://arxiv.org/abs/2305.10429>.
- [36] Mutian Yang, Zisen Zhan, Yutong Chen, Haolin Li, Kaiwen Wang, Kaili Zheng, Yuguang Wang, Qi Wang, Jiandong Gao, and Ji Wu. Learning the mechanism of catastrophic forgetting: A perspective from gradient similarity. *arXiv preprint arXiv:2601.21577*, 2026.
- [37] Cheng Zhang, Hedvig Kjellström, and Stephan Mandt. Determinantal point processes for mini-batch diversification. In *Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence (UAI)*, 2017. URL <https://arxiv.org/abs/1705.00607>.
- [38] Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, Susan Zhang, Gargi Ghosh, Mike Lewis, Luke Zettlemoyer, and Omer Levy. LIMA: Less is more for alignment. In *Advances in Neural Information Processing Systems*, volume 36, pages 55006–55021, 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/hash/ac662d74829e4407ce1d126477f4a03a-Abstract-Conference.html.

Appendices

A Main-body figures	11
B Notation	11
C Experimental details	13
D Taylor analysis of forgetting	15
E Extended LLM results	18
F Compositional capabilities results	23
G Wrappers and concentration	25
H Related Work	31

Appendix A. Main-body figures

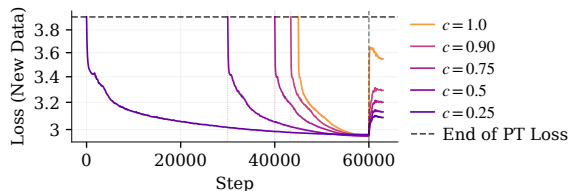


Figure 4: Validation loss during fine-tuning and reversion under a volume-matched budget. Vertical dashed line marks the fine-tuning→reversion boundary. Lower concentration produces monotonically smaller loss rebounds. Biomedical (PubMed) data.

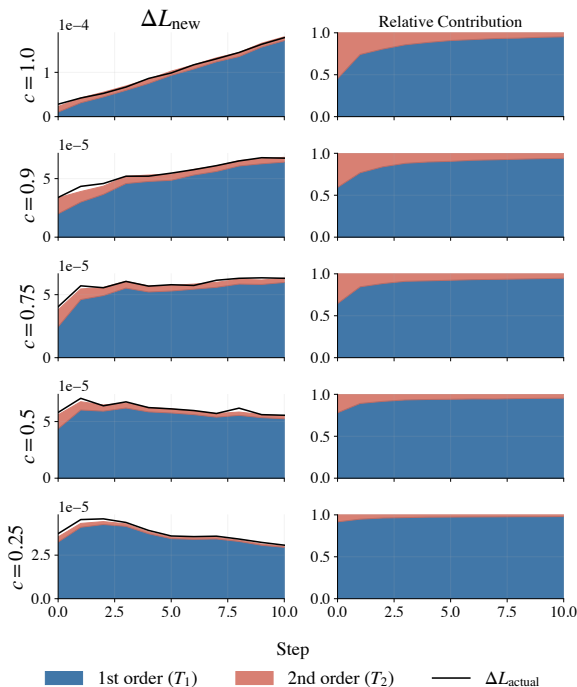


Figure 5: First and second order terms of the Taylor expansion during the first steps of the reversion phase. Higher concentrations lead to higher magnitudes and relative contributions of the curvature term. The black line represents the actual change in loss, and shows the effect of higher order terms. This data is from biomedical (PubMed) data, and is the average of six seeds. For clarity, the scheduler was removed here, and the learning rate was set to its average value during this period.

Appendix B. Notation

This appendix collects every symbol used in the main text and appendices. Symbols are grouped by role.

B.1. Schedule parameters

$c \in [0, 1]$ Concentration. Expected fraction of new-task examples per fine-tuning step (§2). $c=1$ is pure fine-tuning on the new distribution; $c=0$ is pure pre-training.

$\rho \in [0, 1]$ Correlation between the pre-training and new-task distributions used in the CC and PCFG settings (§G). $\rho=1$ gives fully correlated tasks; $\rho=0$ gives independent tasks.

η Learning rate.

B.2. Distributions and losses

$\mathcal{D}_{\text{pre}}, \mathcal{D}_{\text{new}}$ Pre-training and new-task data distributions.

$L_{\text{pre}}(\theta), L_{\text{new}}(\theta)$ Expected losses on \mathcal{D}_{pre} and \mathcal{D}_{new} at parameters θ .

ΔL_{new} Per-step change in the new-task loss during reversion (Eq. 1).

B.3. Parameters and optima

θ, θ_t Model parameters; θ_t denotes parameters at training step t .

$\Delta\theta$ Change in parameters during a reversion training step (on the pre-train data).

θ^* A fine-tune minimum, i.e. a local minimiser of L_{new} (used in the Taylor analysis of §4 and Appendix D).

$\|\theta - \theta^*\|_{H^*}$ Hessian-weighted distance from θ to θ^* , with metric $H^* = \nabla_{\theta}^2 L_{\text{new}}(\theta^*)$.

B.4. Gradient and curvature quantities

$g_{\text{new}} = \nabla_{\theta} L_{\text{new}}(\theta)$ New-task gradient.

$g_{\text{pre}} = \nabla_{\theta} L_{\text{pre}}(\theta)$ Pre-training gradient.

$\cos(g_{\text{new}}, \Delta\theta)$ First-order alignment: cosine similarity between new-task gradient and the reversion weight update.

$H_{\text{new}} = \nabla_{\theta}^2 L_{\text{new}}(\theta)$ Hessian of the new-task loss at θ .

$\lambda_{\text{dir}} = \hat{\Delta}\theta^{\top} H_{\text{new}} \hat{\Delta}\theta$ Directional curvature of L_{new} along the unit reversion step direction $\hat{\Delta}\theta = \Delta\theta / \|\Delta\theta\|$ (the second-order alignment signal). Computed with a finite-difference Hessian-vector product (Appendix D.2).

$\lambda_{\text{avg}} = \text{Tr}(H_{\text{new}}) / d$ Average curvature across all d parameter directions; a reference scale for λ_{dir} .

$\text{Tr}(H_{\text{new}})$ Trace of the new-task Hessian, estimated with the Hutchinson stochastic trace estimator (Appendix D).

B.5. Forgetting and retention metrics

Reversion AUC Trapezoidal area under the new-task accuracy curve during reversion (Appendix C.5). Higher AUC indicates more durable learning.

Life threshold α For $\alpha \in \{0.95, 0.90, \dots, 0.70\}$, the first reversion step at which new-task accuracy drops below $\alpha \cdot a_{\text{peak}}$, where a_{peak} is the peak new-task accuracy at the end of fine-tuning.

Normalised perplexity, NormPPL(t) Rescales each condition to a common $[0, 1]$ range (1.0 = full domain specialisation retained, 0.0 = fully reverted to pre-training baseline).

Appendix C. Experimental details

C.1. Pythia architecture and training schedule

Both Pythia models are decoder-only transformers from the Pythia suite [6], pre-trained on the deduplicated Pile ($\sim 207\text{B}$ tokens, ~ 1.5 epochs) with rotary positional embeddings, parallelised attention and feedforward layers, Flash Attention, and untied input/output embeddings. Table 1 summarises the architecture; Table 2 gives the training schedule across all three phases.

All LLM runs share: Adam ($\beta = (0.9, 0.95)$), $\epsilon = 10^{-8}$), weight decay 0.01, gradient clipping 1.0, and fine-tuning sequence length 512.

	Pythia-70M	Pythia-1B
Total params	70.4M	1011.8M
Non-embedding	18.9M	805.7M
Layers	6	16
Model dimension	512	2048
Attention heads	8	8
Peak LR	10^{-3}	3×10^{-4}
Sequence length	2048	2048

Table 1: Architecture and pre-training hyperparameters for the Pythia models used in our experiments [6]. Both use the deduplicated variant.

Volume-matched control (Pythia-1B). Because low-concentration schedules in the fixed-steps regime see fewer domain examples total, we additionally run a volume-matched experiment where $\text{ft_steps} = 3000/c$ (so all schedules see the same total domain examples): 100% gets 3,000 steps, 50% gets 6,000, 25% gets 12,000.

	Pre-train	Fine-tune	Reversion
<i>Pythia-1B — SMILES</i>			
Steps	143k (original)	10,000	3,000
Peak LR	3×10^{-4}	5×10^{-5}	5×10^{-5}
LR schedule	cosine \rightarrow 10%	constant	constant
Warmup	1% (1,430 steps)	200 steps	—
Eff. batch	1024	16	16
<i>Pythia-70M — SMILES / PubMed</i>			
Steps	143k (original)	12k / 8k	3,000
Peak LR	10^{-3}	5×10^{-5}	5×10^{-5}
LR schedule	cosine \rightarrow 10%	constant	constant
Warmup	1% (1,430 steps)	400 / 200 steps	—
Eff. batch	1024	20 / 40	16 / 40

Table 2: Training schedule for each Pythia experiment. “Pre-train” refers to the original Pythia pre-training [6]; we load the final checkpoint and apply our fine-tune and reversion phases. Where two values are shown (X / Y), the first is SMILES and the second PubMed. Effective batch size accounts for gradient accumulation.

C.2. LLM datasets and budget modes

Models. EleutherAI/pythia-70m-deduped and EleutherAI/pythia-1b-deduped.

Datasets (HuggingFace IDs). Chemistry: antoinebcx/smiles-molecules-chembl. Music: sander-wood/irishman. Biomedical: ccdv/pubmed-summarization. General: monology/pile-uncopyr. All text tokenized with the Pythia tokenizer, chunked into non-overlapping 512-token blocks. Fine-tune mixing is applied *per-sample within each batch*.

Budget modes. In *steps* mode every fine-tuning run uses the same number of steps, so low- c runs see proportionally less new-task data. In *volume* mode the step count scales as $\lceil \text{ft_steps}/b \rceil$ so every fine-tuning run sees the same number of new-task tokens. The 70M sweeps used *steps*; the 1B sweeps used *volume*.

C.3. CC training schedule

Table 3 gives the CC training schedule. All phases use AdamW ($\beta = (0.9, 0.9)$), weight decay 10^{-3} , gradient clipping at 1.0, mixed precision (bfloat16), and batch size 128. The LR follows a piecewise cosine schedule that decays continuously across all three phases. Schedules below $c=0.40$ are excluded from main results because they do not reliably converge to 100% new-task accuracy within the allocated fine-tune phase.

	Pre-train	Fine-tune	Reversion
Steps	600	50	500
Peak LR	10^{-3}	(continuing)	(continuing)
LR at phase end	30% of peak	15% of peak	5% of peak
LR schedule	piecewise cosine (continuous across phases)		
Warmup	50 steps	—	—
Batch size	128	128	128
Seeds	1 (shared)	10 per c	10 per c

Table 3: CC training schedule across all three phases. The LR follows a single piecewise cosine that decays continuously through pre-train, fine-tune, and reversion.

C.4. Concentration and correlation sweep grid values

LLM concentration levels. Pythia-1B (SMILES): $c \in \{1.0, 0.95, 0.90, 0.75, 0.50, 0.25\}$.

Pythia-70M (SMILES): same levels. Pythia-70M (PubMed): $c \in \{1.0, 0.98, 0.95, 0.90, 0.75, 0.50, 0.25, 0.10\}$.

CC concentration levels. $c \in \{1.00, 0.98, 0.95, 0.90, 0.80, 0.70, 0.60, 0.50, 0.40\}$.

CC correlation \times concentration grid. Correlation $\rho \in \{0/n_a, 1/n_a, \dots, n_a/n_a\}$ crossed with $c \in \{100\%, 95\%, 90\%, 70\%, 50\%, 30\%\}$ at $n_a \in \{3, 4, 5, 6\}$ (3 – 5 seeds per cell).

PCFG correlation \times concentration grid.

$$\rho \in \{0.00, 0.25, 0.50, 0.66, 0.75, 0.85, 0.92, 0.95, 1.00\}$$

crossed with

$$c \in \{0.10, 0.30, 0.40, 0.50, 0.60, 0.70, 0.80, 0.85, 0.90, 1.00\}.$$

C.5. Reversion AUC

We measure how quickly the fine-tuned capability is forgotten using *Reversion AUC*, the trapezoidal area under the new-task accuracy curve during reversion. Higher AUC indicates more durable learning.

Appendix D. Taylor analysis of forgetting

D.1. Forgetting step impact on fine-tuning loss

During the continued pre-training stage we perform updates, $\Delta\theta$, to the model using pre-training data. The resulting change in fine-tuning loss can be obtained by Taylor expanding L_{new} around the current weights θ :

$$L_{\text{new}}(\theta + \Delta\theta) = L_{\text{new}}(\theta) + \nabla_{\theta}L_{\text{new}}(\theta)^{\top} \Delta\theta + \frac{1}{2}\Delta\theta^{\top} H_{\text{new}} \Delta\theta + \mathcal{O}(\|\Delta\theta\|^3), \quad (2)$$

where $H_{\text{new}} = \nabla_{\theta}^2 L_{\text{new}}(\theta)$ is the Hessian of the fine-tuning loss and $g_{\text{new}} = \nabla_{\theta}L_{\text{new}}(\theta)$. We rearrange, and identify the first and second order terms,

$$\Delta L_{\text{new}} = \underbrace{g_{\text{new}}^{\top} \Delta\theta}_{T_1} + \underbrace{\frac{1}{2}\Delta\theta^{\top} H_{\text{new}} \Delta\theta}_{T_2} + \mathcal{O}(\|\Delta\theta\|^3). \quad (3)$$

The first term, T_1 , is the dot product between the new-task gradient and the weight update: it is negative when the weight update opposes the new-task gradient, and positive when they align. The second term, T_2 , captures how curved the fine-tuning loss landscape is in the direction the pre-training step takes the model. At a minima, this term is always non-negative, meaning curvature always contributes to forgetting regardless of whether the gradients are aligned.

D.2. Relative magnitudes of T_1 and T_2 at a fine-tuning minimum

We now show that, near a fine-tuning minimum, T_1 is bounded by the distance from that minimum — implying that models fine-tuned at higher concentration, which sit closer to the minimum, will have forgetting dominated by the curvature term T_2 .

To consider the relative magnitudes of the first and second order effects, suppose $c = 1$, and assume we will reach a minima in the new-task loss, L^* , with parameters, θ^* . As one moves away from a pure concentration of fine-tuning data, $c < 1$, the model is also be optimizing for the pre-training data, and thus at the end of fine-tuning, it may not be so deep in the fine-tuning minima, and there will be a non-zero displacement term between its weights, θ , and the fine-tuning minima, i.e. $(\theta - \theta^*) \neq 0$. Finding the Taylor expansion of the fine-tuning loss around this fine-tuned minimum yields a primary term that is quadratic,

$$L_{\text{new}}(\theta) \approx L^* + \frac{1}{2}(\theta - \theta^*)^{\top} H^*(\theta - \theta^*) + \mathcal{O}(\|(\theta - \theta^*)\|^3). \quad (4)$$

Thus, ignoring higher order terms, the gradient of the loss in this region is approximately linear in the displacement from the minima,

$$\frac{\partial L_{\text{new}}(\theta)}{\partial \theta} = g_{\text{new}}(\theta) \approx H^*(\theta - \theta^*). \quad (5)$$

Now considering a forgetting step, $\Delta\theta$ at the current parameters, θ , the first and second order contributions to ΔL_{ft} are,

$$T_1 = g_{\text{new}}^{\top} \Delta\theta \approx (\theta - \theta^*)^{\top} H^* \Delta\theta, \quad T_2 = \frac{1}{2}\Delta\theta^{\top} H^* \Delta\theta. \quad (6)$$

Notice the asymmetry: T_1 depends on the displacement $(\theta - \theta^*)$ but T_2 does not.

Applying the Cauchy-Schwarz Inequality¹, one obtains a relation between the magnitudes of these terms,

$$|T_1| \leq \sqrt{2T_2} \cdot \|\theta - \theta^*\|_{H^*}. \quad (7)$$

This bounds the first-order term by the product of the curvature term and the distance from the minimum. We note that this is general for any optimizer used (Adam, SGD, etc.). The implication is:

At *high concentration* ($c \rightarrow 1$), the model ends fine-tuning close to θ^* , so $\|\theta - \theta^*\|_{H^*}$ is small. The bound is tight: T_1 is forced to be small relative to T_2 , and forgetting is driven primarily by curvature. At *low concentration*, the model is pulled away from θ^* by the pre-training objective, so $\|\theta - \theta^*\|_{H^*}$ is large. The bound loosens: T_1 can contribute substantially, and forgetting reflects both gradient conflict and curvature.

D.3. Calculations: Taylor-term and curvature estimation

For each end-of-fine-tuning checkpoint $\theta_{\text{ft}}^{(c)}$, we simulate N_{steps} steps of the reversion phase on the pre-train loss starting from $\theta_{\text{ft}}^{(c)}$. At each step $k \in \{0, \dots, N_{\text{steps}}\}$ we record the two terms of the Taylor expansion of L_{ft} around the current iterate θ_k :

$$T_1(\theta_k) = \langle g_{\text{new}}(\theta_k), \Delta\theta_k \rangle, \quad (8)$$

$$T_2(\theta_k) = \frac{1}{2} \|\Delta\theta_k\|^2 \lambda_{\text{dir}}(\theta_k), \quad (9)$$

where g_{ft} is the gradient of the fine-tune loss, and $\lambda_{\text{dir}}(\theta_k) = \widehat{\Delta\theta}^\top H_{\text{ft}}(\theta_k) \widehat{\Delta\theta}$ is the Rayleigh quotient of the fine-tune Hessian along the unit pre-train gradient direction, $\widehat{\Delta\theta} = \Delta\theta / \|\Delta\theta\|$. After recording, we take one optimizer step (using Adam) along the pre-train gradient,

$$\theta_{k+1} = \theta_k + \eta \cdot \text{Adam}(L_{\text{pre}}(\theta_k)),$$

and repeat.

Gradient estimation. For the proceeding sections, due to the learning rate scheduler, the initial learning rate is 0, and so to calculate $\Delta\theta$, we instead use a value that is the average of the first 10 steps of the scheduler.

Both g_{new} and $\Delta\theta$ are estimated as the mean of the per-batch gradients over N_{bat} mini-batches of size B , drawn from the held-out validation splits of the domain (fine-tune) and pile (pre-train) datasets respectively.

Directional curvature λ_{dir} . Rather than instantiating H_{ft} , we compute Hessian-vector products with a central finite-difference approximation:

$$H_{\text{ft}} v \approx \frac{g_{\text{ft}}(\theta + \varepsilon v) - g_{\text{ft}}(\theta - \varepsilon v)}{2\varepsilon}, \quad \|v\| = 1,$$

with $\varepsilon = \text{FD_EPS}$. The parameter vector is restored to θ after each probe. Setting $v = \widehat{\Delta\theta}$ gives $\lambda_{\text{dir}} = \widehat{\Delta\theta}^\top (H_{\text{ft}} \widehat{\Delta\theta})$.

Average curvature λ_{avg} . For the post-fine-tune checkpoint of each concentration, we additionally estimate the average eigenvalue of H_{ft} via Hutchinson’s trace estimator. We

1. The Cauchy-Schwarz inequality $|\langle a, b \rangle| \leq \|a\| \|b\|$ holds for any inner product. Taking $a = \theta - \theta^*$, $b = \Delta\theta$, and the inner product $\langle a, b \rangle_{H^*} = a^\top H^* b$:

$$\begin{aligned} |T_1| &= |\langle \theta - \theta^*, \Delta\theta \rangle_{H^*}| \leq \|\theta - \theta^*\|_{H^*} \|\Delta\theta\|_{H^*} \\ &= \|\theta - \theta^*\|_{H^*} \sqrt{2T_2}, \end{aligned}$$

where we used $\|\Delta\theta\|_{H^*}^2 = \Delta\theta^\top H^* \Delta\theta = 2T_2$.

draw N_{hutch} Rademacher vectors $v^{(i)} \in \{-1, +1\}^d$ ($d = \text{total parameter count}$), form unit vectors $\hat{v}^{(i)} = v^{(i)} / \|v^{(i)}\| = v^{(i)} / \sqrt{d}$, and compute the same finite-difference Rayleigh quotient $\kappa^{(i)} = (\hat{v}^{(i)})^\top H_{\text{ft}} \hat{v}^{(i)}$ for each. Because $\mathbb{E}_{\hat{v}}[\hat{v}^\top H \hat{v}] = \text{Tr}(H)/d$ for isotropic unit directions, we estimate

$$\lambda_{\text{avg}} = \frac{\text{Tr}(H_{\text{ft}})}{d} \approx \frac{1}{N_{\text{hutch}}} \sum_{i=1}^{N_{\text{hutch}}} \kappa^{(i)}.$$

The full trace estimate is recovered as $d \cdot \lambda_{\text{avg}}$.

Hyperparameters. We use the same reversion learning rate as the main experiments, $N_{\text{steps}} = 10$, $B = 16$, $N_{\text{bat}} = 10$, $\varepsilon = 10^{-3}$, and $N_{\text{hutch}} = 50$. Further, we use 6 seeds for the calculations of λ_{dir} and λ_{avg} , and the error bars on the plots represent ± 1 standard deviation. All computations are performed in single-precision float32.

D.4. Update magnitude under SGD vs. Adam

Under vanilla SGD, the reversion update is $\Delta\theta = -\eta g_{\text{pre}}$, so $\|\Delta\theta\| \propto \|g_{\text{pre}}\|$ exactly and any pre-training drift accumulated during fine-tuning translates directly into larger reversion steps. For adaptive optimisers such as Adam, this coupling is weaker. Adam’s update is approximately $\Delta\theta \approx -\eta \hat{m} / (\sqrt{\hat{v}} + \epsilon)$, where \hat{m} and \hat{v} are bias-corrected first- and second-moment estimates of the gradient. At the first reversion step, the bias correction makes the update approximately a per-coordinate sign step, decoupling $\|\Delta\theta\|$ from $\|g_{\text{pre}}\|$. As reversion proceeds, \hat{v} accumulates and begins to track the pre-training gradient statistics, restoring partial coupling between update magnitude and gradient norm.

This means $\|g_{\text{pre}}\|$ is a clean drift proxy under SGD but only a qualitative one under Adam, particularly in the first few reversion steps. Despite this, we observe empirically (Fig. 6) that the reversion update norms remain monotonic in c throughout the early reversion window, indicating that the concentration ordering survives Adam’s normalisation. We attribute this to the fact that higher- c fine-tuning increases $\|g_{\text{pre}}\|$ by a large enough margin (Fig. 3b, especially the $c = 1$ outlier) that even Adam’s per-coordinate normalisation does not erase the ordering.

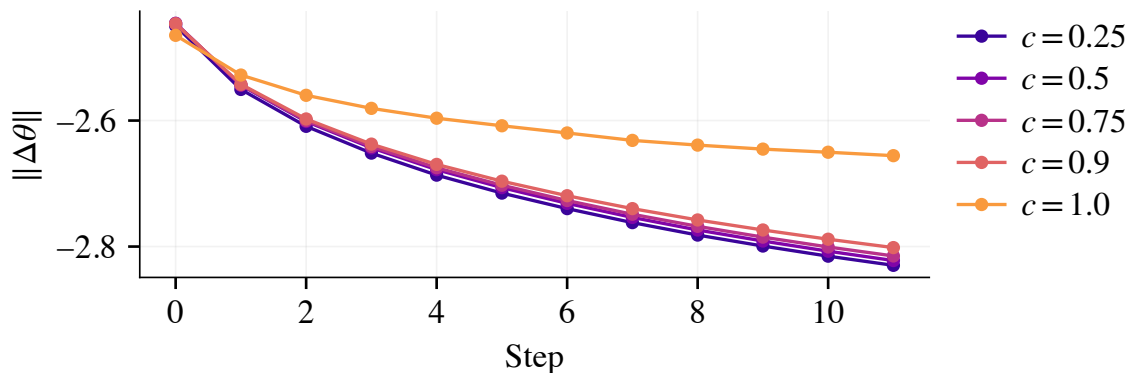


Figure 6: The magnitude of the weight update across the 11 steps of training from Fig. 5. One sees that it is higher for larger concentrations.

We also show additional results (Fig. 7) for the relative magnitudes of the T_1 and T_2 when SGD is used instead of AdamW, showing the same trend as Fig. 5.

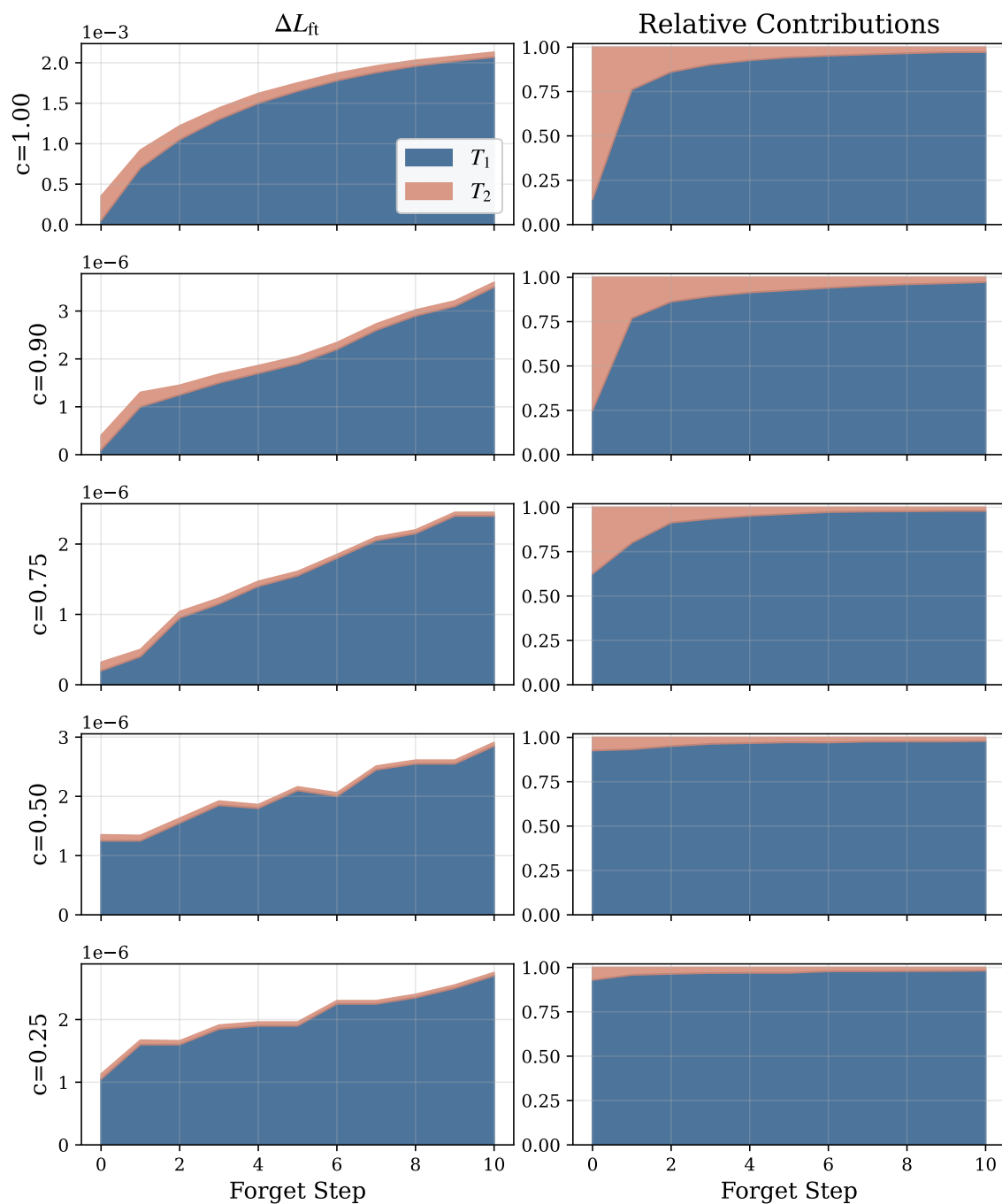


Figure 7: Analogous figure to Fig. 5, but for the case of Vanilla SGD instead of AdamW.

Appendix E. Extended LLM results

This appendix collects extended LLM results that complement the main findings in Section 3. Unless specified otherwise, all the subsequent results are for Pythia-70M fine-tuned on Chemistry (SMILES, ChEMBL), $\text{lr} = 5 \times 10^{-5}$, fine-tuning for 15k steps, and continued pre-training for 3k steps.

c	Chemistry	Music	Biomedical
	FT=15k, CPT=3k	FT=15k, CPT=3k	FT=15k, CPT=3k
1.00	0.20	2.43	1.45
0.90	0.12	1.59	0.83
0.75	0.09	1.25	0.64
0.50	0.07	0.87	0.47
0.25	0.05	0.62	0.36

Table 4: Relative degradation: $\ell_{\text{CPT}}/\ell_{\text{new}} - 1$, where ℓ_{new} is the domain validation loss at the end of fine-tuning and ℓ_{CPT} at the end of continued pre-training. Higher values indicate more degradation.

c	Post-FT		Post-CPT	
	Domain	Pile	Domain	Pile
Pretrained	3.92	3.06	—	—
1.00	2.96	5.92	3.55	3.02
0.90	2.97	3.10	3.31	3.01
0.75	2.98	3.04	3.25	3.00
0.50	3.02	3.01	3.23	2.99
0.25	3.08	2.99	3.25	2.99

Table 5: Validation cross-entropy loss after fine-tuning (Post-FT) and after continued pre-training on the Pile (Post-CPT). Lower concentration keeps lower loss post continued pre-training.

Volume-matched chemistry control (Pythia-1B). The main text shows the volume-matched PubMed control in Fig. 4. For completeness, we include the analogous Pythia-1B chemistry (SMILES) run here. Matching the total number of domain tokens across concentrations again preserves a monotonic relationship between concentration and the post-reversion rebound in domain loss.

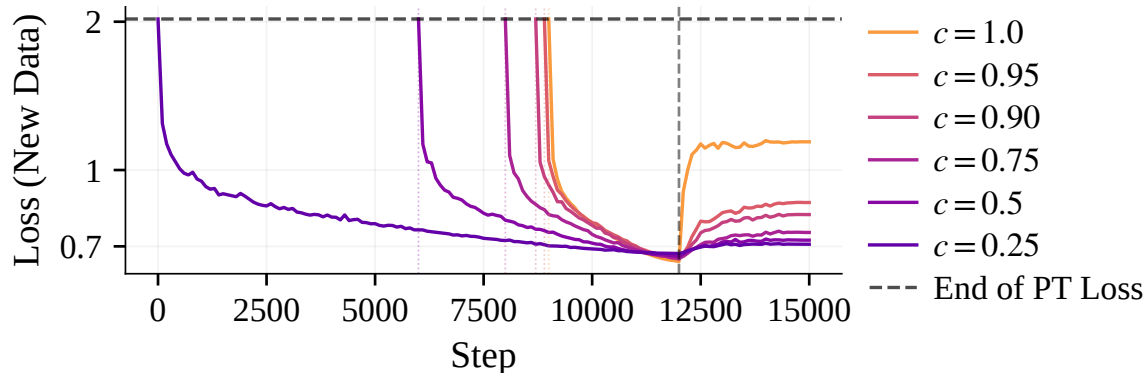


Figure 8: Pythia-1B domain validation loss on Chemistry (SMILES) during fine-tuning and continued pre-training on a volume-matched budget. As in the PubMed control shown in the main text, lower concentration yields smaller loss rebounds after the fine-tune→reversion boundary.

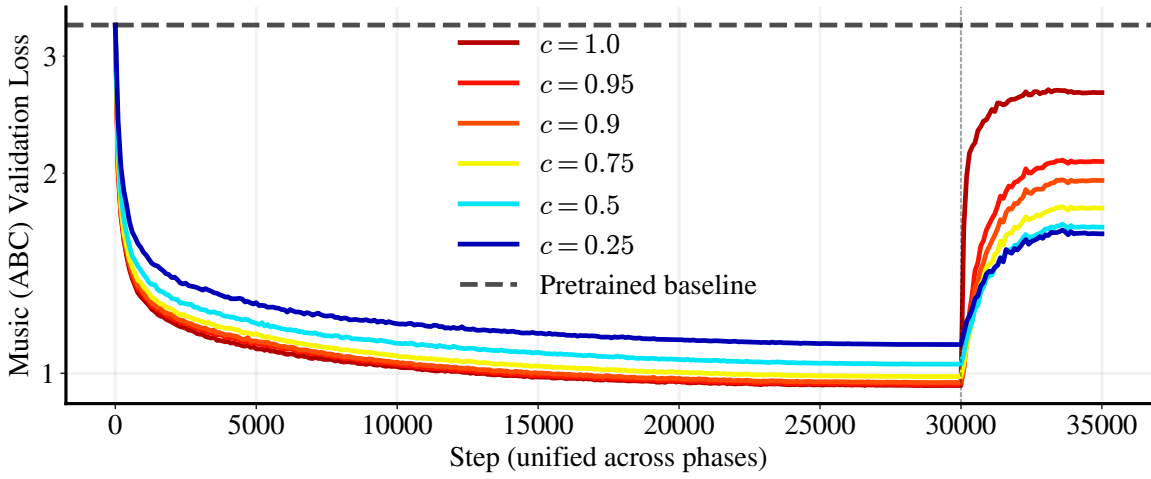


Figure 9: Domain validation loss, fine-tuning on Music (ABC notation).

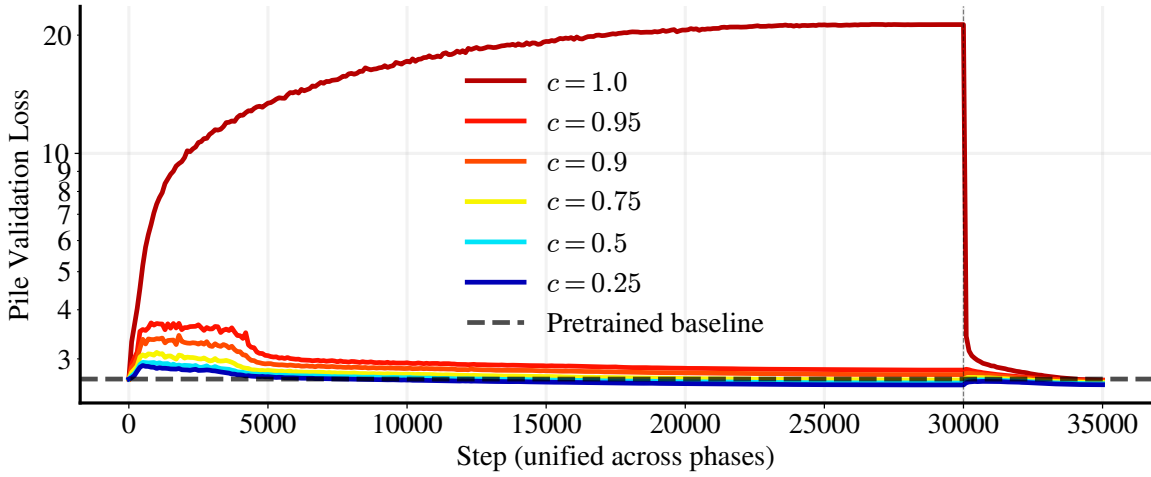


Figure 10: Pre-training data (Pile) Validation loss, fine-tuning. The pre-training loss is disproportionately higher in the 100% case.

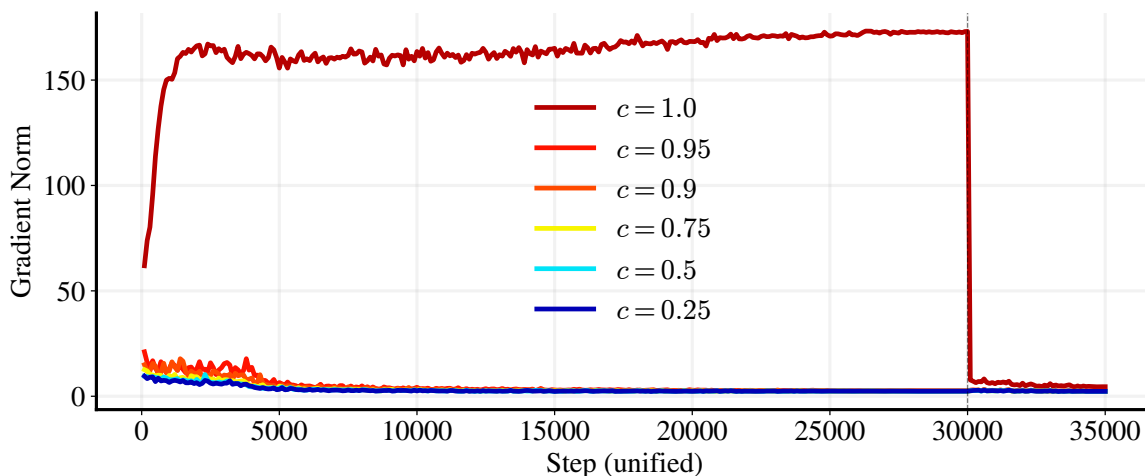


Figure 11: Gradient norm of the pre-training data, fine-tuning. The pre-training gradient norm is disproportionately higher in the 100% case.

Stepwise PubMed diagnostic. The mechanistic analysis in the main text uses Pythia-70M fine-tuned on PubMed to show that the pre-training task gradient norm grows with concentration. We include the full stepwise trace here for reference.

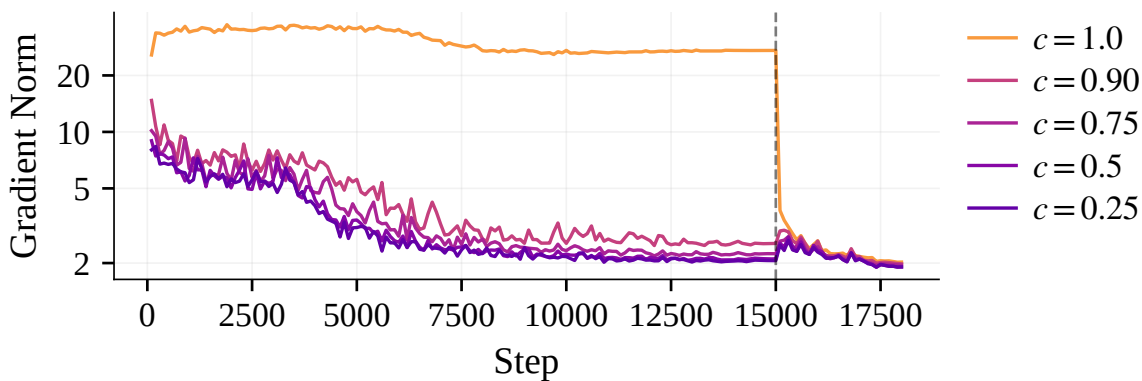


Figure 12: Biomedical dataset. Pre-training task gradient norm during fine-tuning.

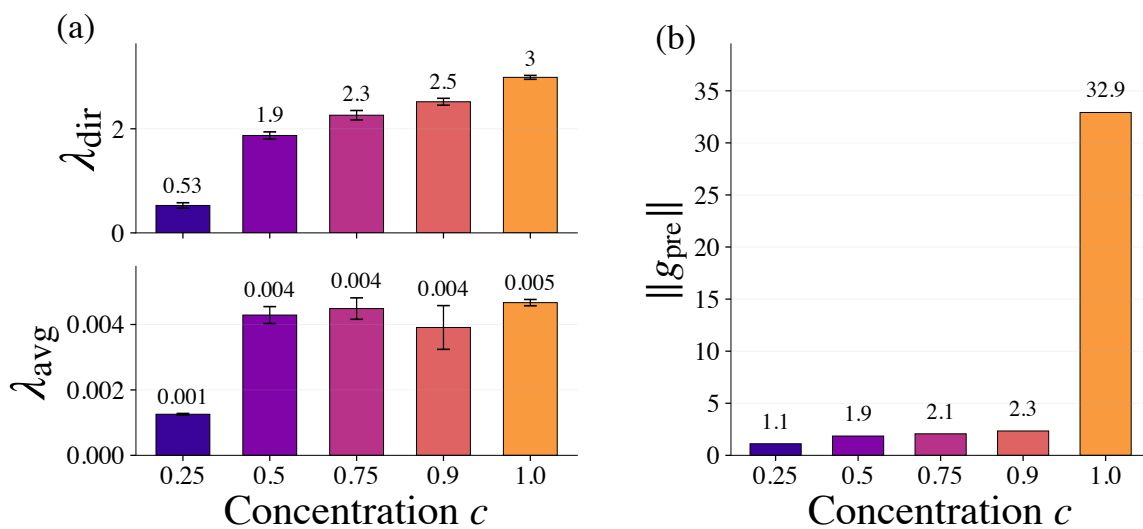


Figure 14: The same plot as in Fig. 3 in the main text, but (a) is shown in linear scale. The error bars are the $\pm 1\sigma$ from 6 seeds.

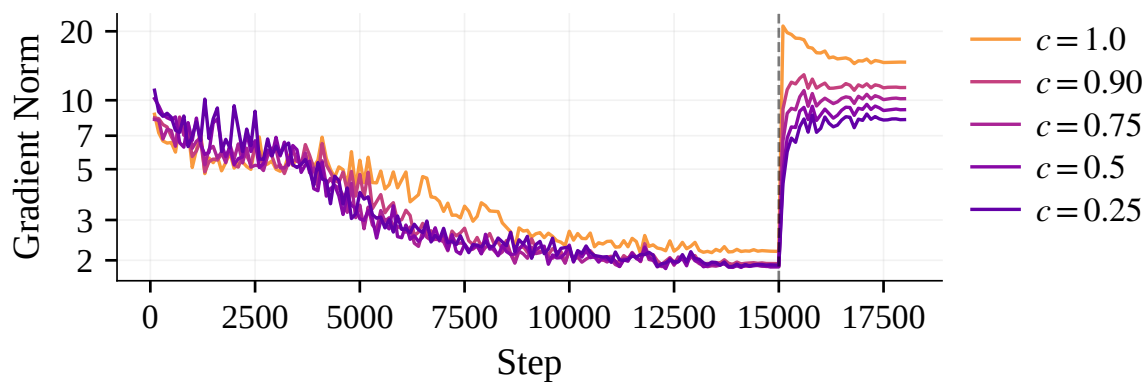


Figure 13: Biomedical dataset. Pre-training task gradient norm during PubMed fine-tuning, shown across training steps for each concentration. Higher concentration produces larger pre-training gradient norms throughout fine-tuning, matching the trend summarized in Fig. 12.

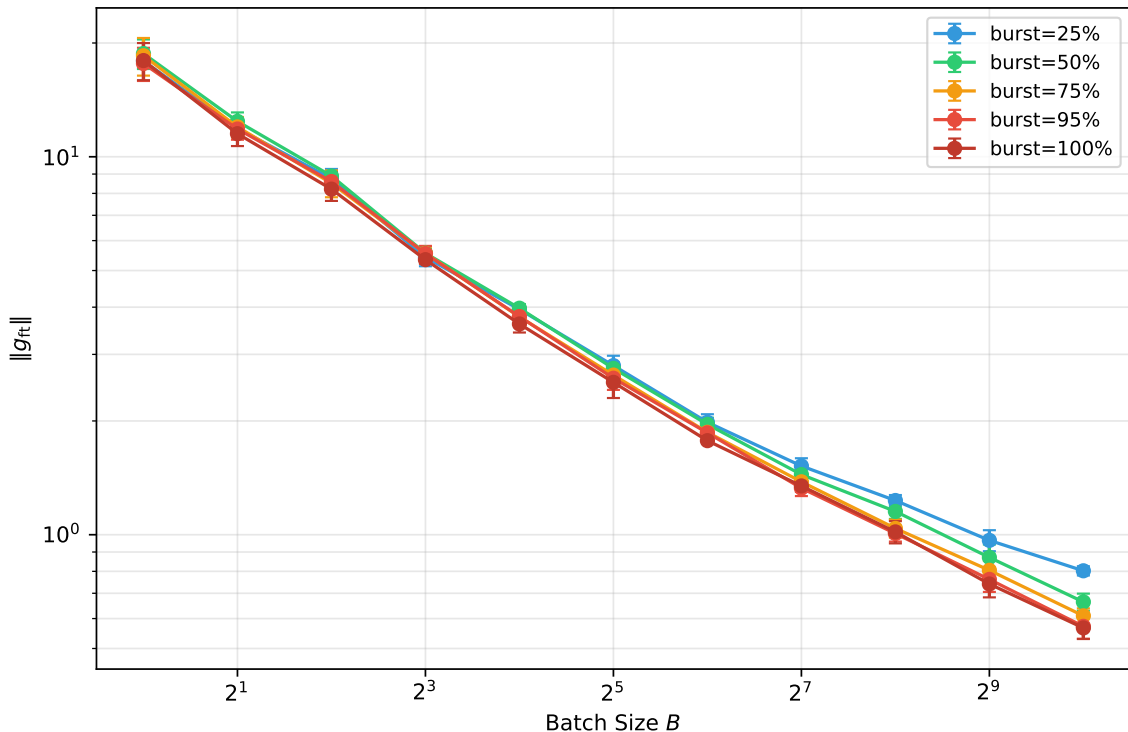


Figure 15: Gradient norm of the new-task data after fine-tuning. One sees the norm follows a \sqrt{B} relation, and when the data is sufficiently represented, the norm of higher concentrations is lower than that of the lower concentrations.

E.1. Learning rate sensitivity

The magnitude of forgetting depends strongly on hyperparameter choices, especially learning rate [13, 16] and optimizer [27]. We also empirically observed that batch size had a significant effect. To quantify this interaction, we sweep the learning rate across the fine-tuning phase on the Compositional Capabilities setup while holding all other hyperparameters fixed, and measure forgetting at several concentration levels (Fig. 16).

The learning rate has a dramatic effect on forgetting behaviour: with this setting, typically at low learning rates (relative to pre-training learning rate) fine-tuning is largely forgotten during reversion, while at high learning rates the fine-tuned capability is retained almost entirely for the allocated steps.

Our monotonic ordering replicates across the settings we tested, but absolute numbers should not be compared across hyperparameter regimes.

Appendix F. Compositional capabilities results

This appendix presents the full Compositional Capabilities (CC) experimental setup and results, which complement the LLM findings in the main text. All figures use 10 seeds with concentration schedules $c \in \{40, 50, 60, 70, 80, 90, 95, 98, 100\}\%$ unless noted otherwise.

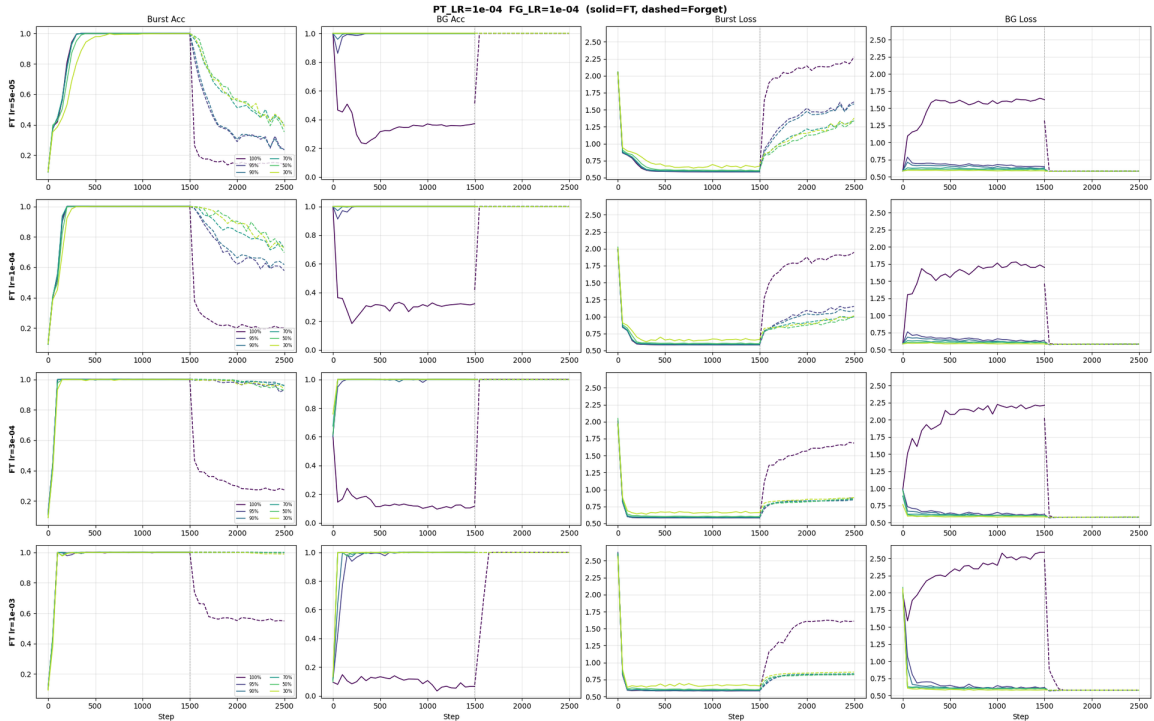


Figure 16: Learning Rate Sweep on Compositional Capabilities showing the effect of this parameter on forgetting.

F.1. Setup

We use the Compositional Capabilities framework of Ramesh et al. [31], in which 6-layer nanoGPT transformers are trained from scratch on pre-training bijection chains. During fine-tuning, a novel bijection is injected into the training mixture at concentration c , with the remaining $1-c$ fraction drawn from the original pre-training distribution. After fine-tuning, training reverts to the pre-training distribution and we measure how quickly the new task is forgotten. Training details (optimizer, learning rate schedule, batch size) are given in Table 3.

F.2. Concentration drives forgetting

All concentration schedules learn the new task to near-perfect accuracy (Table 6), but they differ sharply in how long that knowledge survives reversion. Fig. 17 shows both the Reversion AUC and the new-task accuracy trajectories during the early reversion window: higher concentration produces faster forgetting across all metrics. The effect is graded and monotonic — even small reductions in c (e.g. from 100% to 95%) yield measurably more durable learning. Pre-training task accuracy remains at 1.0 throughout for all schedules, confirming that lower concentrations do not impair the original task.

F.3. Mechanistic diagnostics

We report the same gradient diagnostics as in the main text, applied to the CC setting.

Gradient alignment. Consistent with Yang et al. [36], the cosine similarity between new and pre-training gradients becomes increasingly negative during fine-tuning (Fig. 18), indicating

Schedule	Peak	Rev. AUC	90%-life	80%-life	70%-life
Concentration 100%	1.000	216 \pm 5	25 \pm 0	25 \pm 0	25 \pm 0
Concentration 98%	1.000	254 \pm 8	25 \pm 0	25 \pm 0	32 \pm 7
Concentration 95%	1.000	286 \pm 4	25 \pm 0	40 \pm 8	60 \pm 8
Concentration 90%	1.000	328 \pm 8	40 \pm 8	80 \pm 14	150 \pm 21
Concentration 80%	1.000	362 \pm 7	72 \pm 15	148 \pm 16	238 \pm 31
Concentration 70%	1.000	374 \pm 7	78 \pm 11	150 \pm 14	272 \pm 40
Concentration 60%	1.000	376 \pm 4	72 \pm 15	172 \pm 16	280 \pm 23
Concentration 50%	1.000	376 \pm 5	80 \pm 22	165 \pm 10	275 \pm 31
Concentration 40%	1.000	378 \pm 4	85 \pm 10	175 \pm 21	268 \pm 25

Table 6: CC task summary statistics (10 seeds, \pm 95% CI). Peak accuracy is near-perfect for all schedules. Reversion AUC and $X\%$ -life thresholds (steps to decay below $X\%$ of peak) both increase monotonically with decreasing concentration, quantifying the durability advantage of diluted training.

active competition in parameter space. However, at $c = 1$ the cosine remains near zero while performance still degrades, demonstrating that first-order gradient conflict is not the only driver of forgetting.

Pre-training gradient norm. The pre-training gradient norm grows during fine-tuning and is monotonic with concentration (Fig. 18), consistent with the model moving further from the pre-training optimum at higher c . This matters for both Taylor terms of Eq. 1 in the main text, but unevenly: T_1 is linear in $\|g_{\text{pre}}\|$ while T_2 is quadratic, so larger g_{pre} at higher concentration amplifies the second-order contribution disproportionately.

Taylor decomposition. Fig. 19 decomposes the per-step change in fine-tuning loss during the first reversion steps into its first-order (T_1) and second-order (T_2) Taylor terms. As in the LLM setting, the curvature term’s share of the total grows with concentration, providing converging evidence that the balance between first- and second-order forgetting shifts systematically with c . This analysis was done using vanilla SGD.

Appendix G. Wrappers and concentration

This section details additional results on the effect of varying fine-tuning data concentration on wrapper formation and forgetting under varying correlation. These experiments were mostly performed using the setup from [14] (as well as the CC setup), where a transformer is trained to solve tasks involving probabilistic context free grammars (PCFGs). In the original work, it was shown that wrappers form in the presence of spurious correlations between pre-training and fine-tuning tasks. In these experiments, we test how wrapper formation changes in the case where the fine-tuning data concentration is reduced, as well as probing how forgetting changes with simultaneous sweeps through concentration and correlation.

G.1. PCFG experiment details and training schedule

The setup from [14] involves generating PCFG strings consisting of the characters $\in [a, b, c]$, and then training a transformer to perform tasks such as counting occurrences of a , indexing the n^{th} occurrence of a or finding what token is at position x . The first tokens of the input prompt to the transformer are dedicated to setting the task. For example, a full training

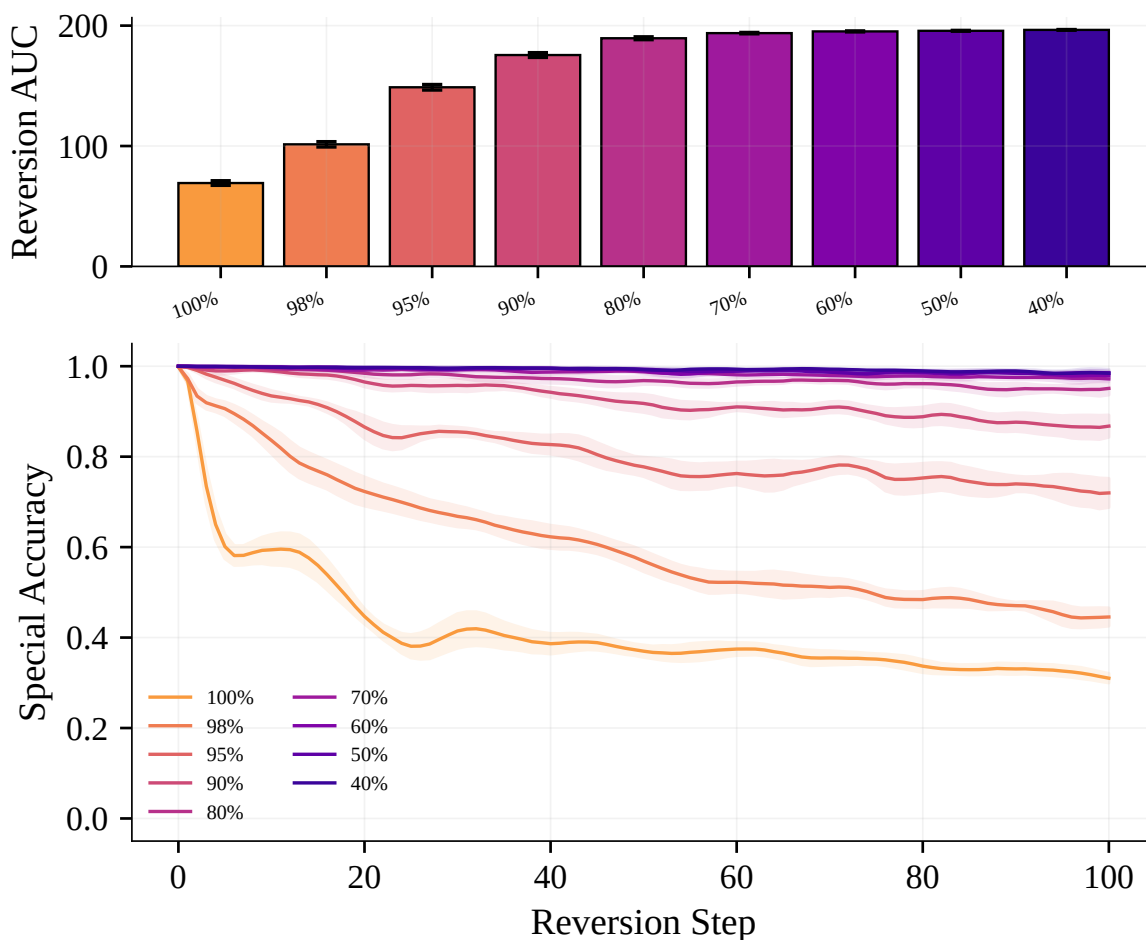


Figure 17: CC summary (10 seeds). **(A)** Reversion AUC across concentrations. **(B)** New-task accuracy during the early reversion window (shaded: 95% CI). Both panels show a monotonic, graded concentration effect: higher concentration yields faster forgetting, and the relationship is smooth rather than a sharp $c=1$ outlier.

string has the form of

$$\$, T, a, x, <, [\text{PCFG string}], >, = y, \$.$$

$\$$ is a start of sequence token, T is a task specific token, a is an operand token, x is a parameter token, and y is the answer token. For example, the string, ‘Ca40’, means to count the occurrences of ‘a’ in the last 40 tokens, ‘Ib6’ means to find the position of the 6th occurrence of ‘b’, and ‘@_40’ means to find the token at the 40th position from the end.

As in the original paper, the PCFG generation length was 255 and the tasks analyzed were:

- Ca40, Cb40, Cc40,
- Caa40, Cbb40, Ccc40
- Ia6, Ib6, Ic6
- Iaa6, Ibb6, Icc6

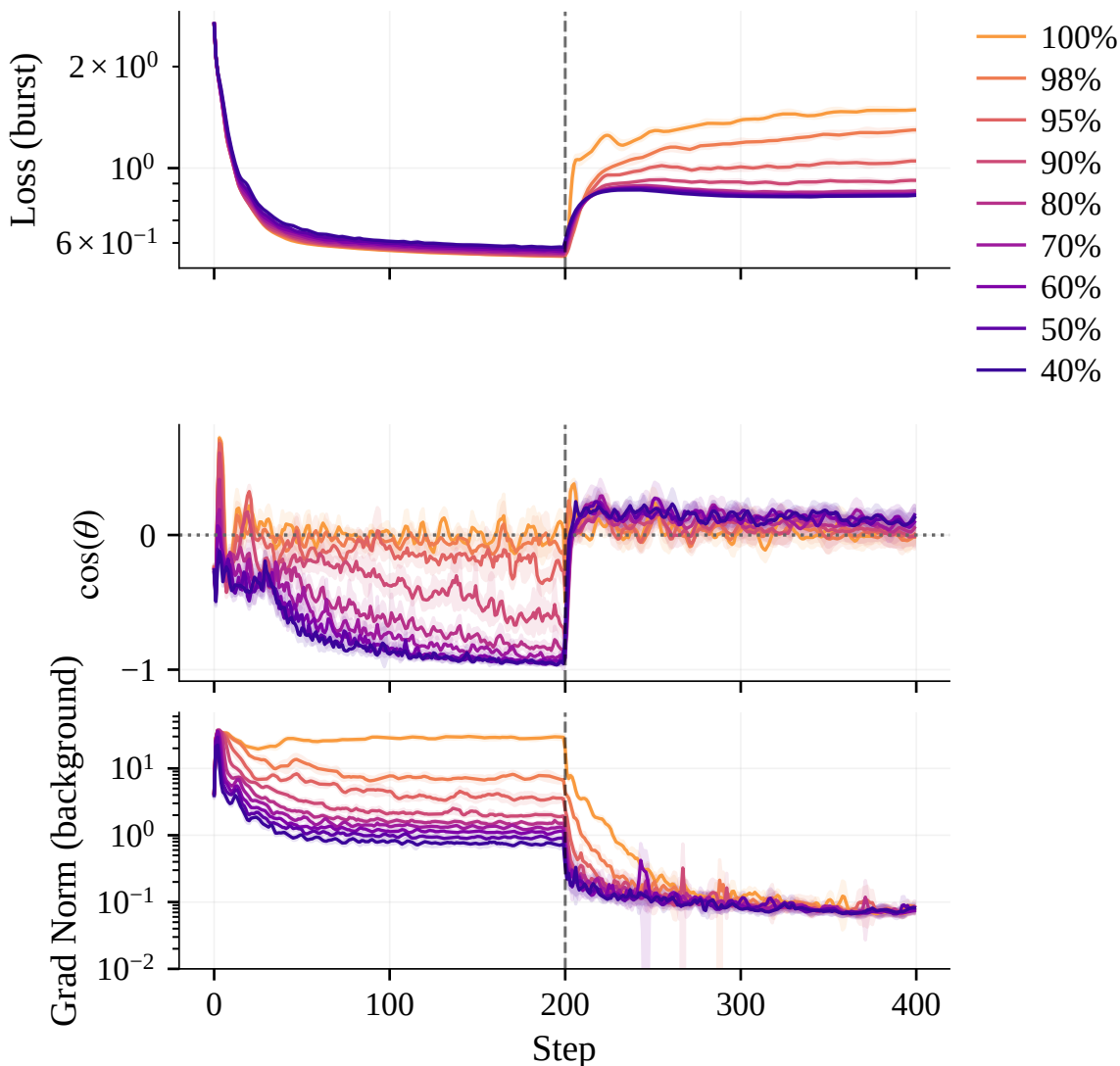


Figure 18: CC gradient diagnostics through fine-tuning and the first 200 reversion steps (log scale), averaged over 10 seeds. Panels show the new-task loss, $\cos(g_{\text{new}}, g_{\text{pre}})$, the pre-training gradient norm $\|g_{\text{pre}}\|$, and the effective rank of the new-task gradient. Higher concentration produces larger pre-training gradient norms and increasingly negative gradient cosine during fine-tuning.

- @_40,

where ‘aa’ refers to an occurrence of consecutive characters. The fine-tune task is chosen to be Ca40, and the pre-train tasks are the remaining 12 tasks. As in the original paper, we implement the correlation, ρ , between the pre-training and fine-tuning data by enforcing that $\text{Ca}40 = \text{Cb}40 + 1$ for a fraction, ρ , of the data. This is done by generating many PCFG strings and extracting ones where the condition is true.

[14] demonstrated the formation of wrappers in the case of $\rho = 1$ by taking the model after fine-tuning, when it had 100% accuracy on Ca40 and 0% accuracy on Cb40, and selectively pruning neurons to improve the performance on Cb40. They demonstrated that pruning just ≈ 10 neurons was sufficient to return the model to its performance at the end of pre-training

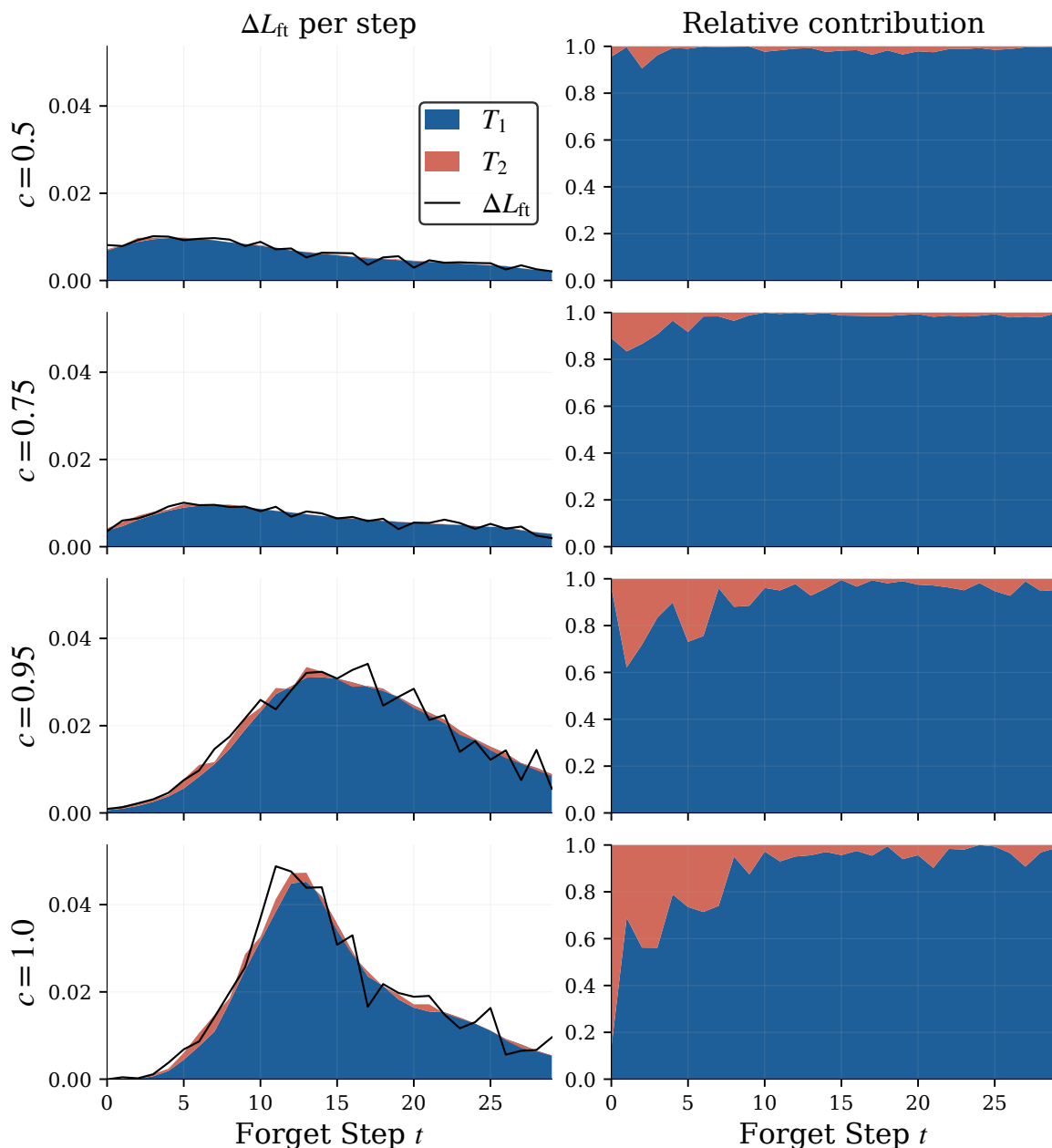


Figure 19: First- and second-order terms of the Taylor expansion during the first reversion steps on the CC task. Higher concentrations produce a larger relative contribution from the curvature term T_2 .

(0% accuracy on Ca40 and 100% accuracy on Cb40). In the results, we extend this analysis by testing wrapper formation in lower concentrations than $c = 1$.

Table 7 gives the PCFG training schedule. All phases use AdamW ($\beta = (0.9, 0.95)$), gradient clipping 1.0), batch size 96, and block size 512.

We also implement correlation in the case of the compositional capabilities setup (CC). Here, for the fine-tuning tasks we introduce six new possible functions. To create correlation,

	Pretrain	Fine-tune	Reverse
Steps	200,000	10,000	8,000
Peak LR	10^{-3}	5×10^{-5}	5×10^{-5}
LR schedule	cosine decay	constant	constant
Batch size	96	96	96
Dropout	0.1	0.1	0.1

Table 7: PCFG training schedule across all three phases.

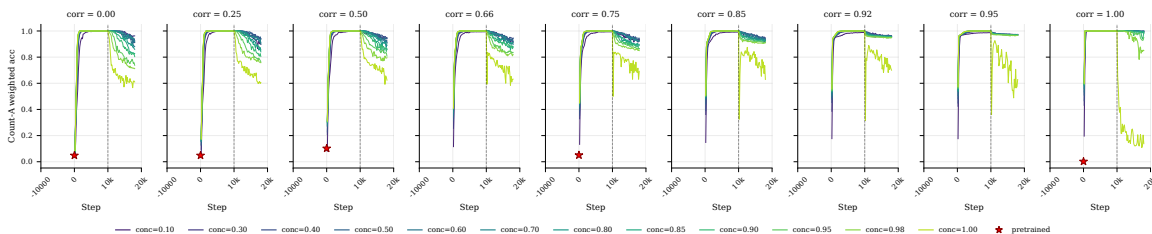


Figure 20: The resulting accuracy on Ca40 and Cb40 after pruning neurons on the PCFG setup with varying values of ρ and c . In the case of $c = 1$, pruning ~ 10 neurons recovers Cb40 accuracy while collapsing Ca40, replicating the wrapper-formation result of [14]. Reducing the concentration seems to avoid the formation of this wrapper.

we simply make a fraction of these functions the same as the pre-training functions, meaning that we obtain 7 possible correlation values.

G.2. Results

Correlation. Firstly, we performed multiple training runs, simultaneously scanning through correlation and concentration on both the PCFG and CC setups; the accuracy on the fine-tune task for the PCFG task is shown in Fig. 20, and the forgetting result for both setups is seen in Fig. 21. We observe what has been discussed in the literature, that typically tasks with higher similarity are forgotten less fast. Further, the results suggest that when training on a task that is more correlated with the pre-training data, one can use a higher value of c than with a task that is less correlated.

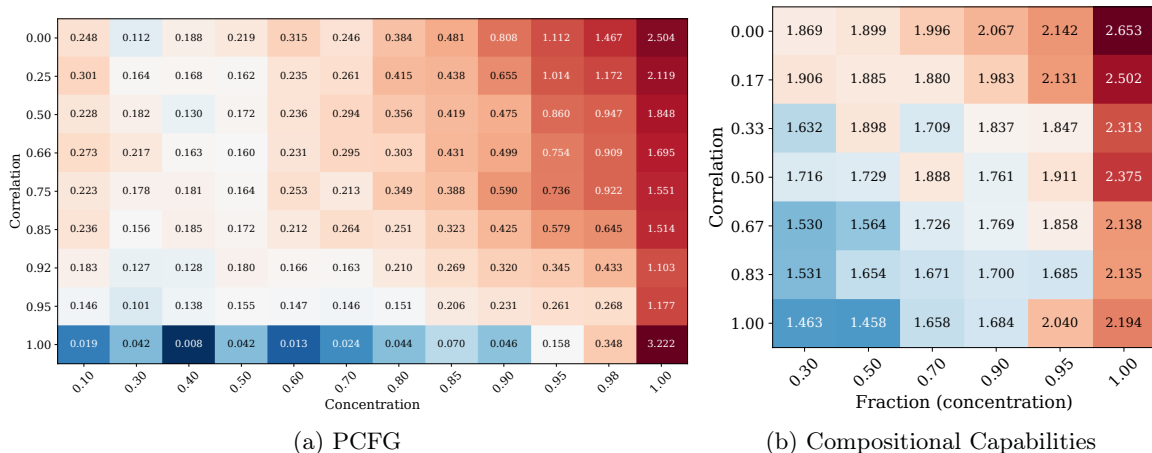


Figure 21: Correlation vs concentration scan. **Left:** PCFG Count-A final reversion loss. **Right:** CC accuracy drop (end vs peak). Both tasks show the same interaction: fragility grows with concentration and shrinks with correlation.

Wrappers. Second, we perform an experiment to investigate wrapper formation at lower concentrations. Firstly, at parameters $[\rho = 1, c = 1]$, we replicate the wrapper-formation result from [14], that pruning 10 neurons from the network causes the model to regain its performance on Cb40, and dramatically to lose its performance on the fine-tune task, Ca40. A similar result is found (to a lesser extent) when reducing ρ to 0.95 — both results are seen in Fig. 22 (for the case of $c = 1$). We then perform the same analysis on models trained with lower concentration $\in [0.9, 0.95]$, and find that this is sufficient to prevent wrapper formation. By the end of fine-tuning, the accuracy on the Cb40 task has not degraded relative to the end of pre-training, and pruning neurons does not remove the performance on Ca40.

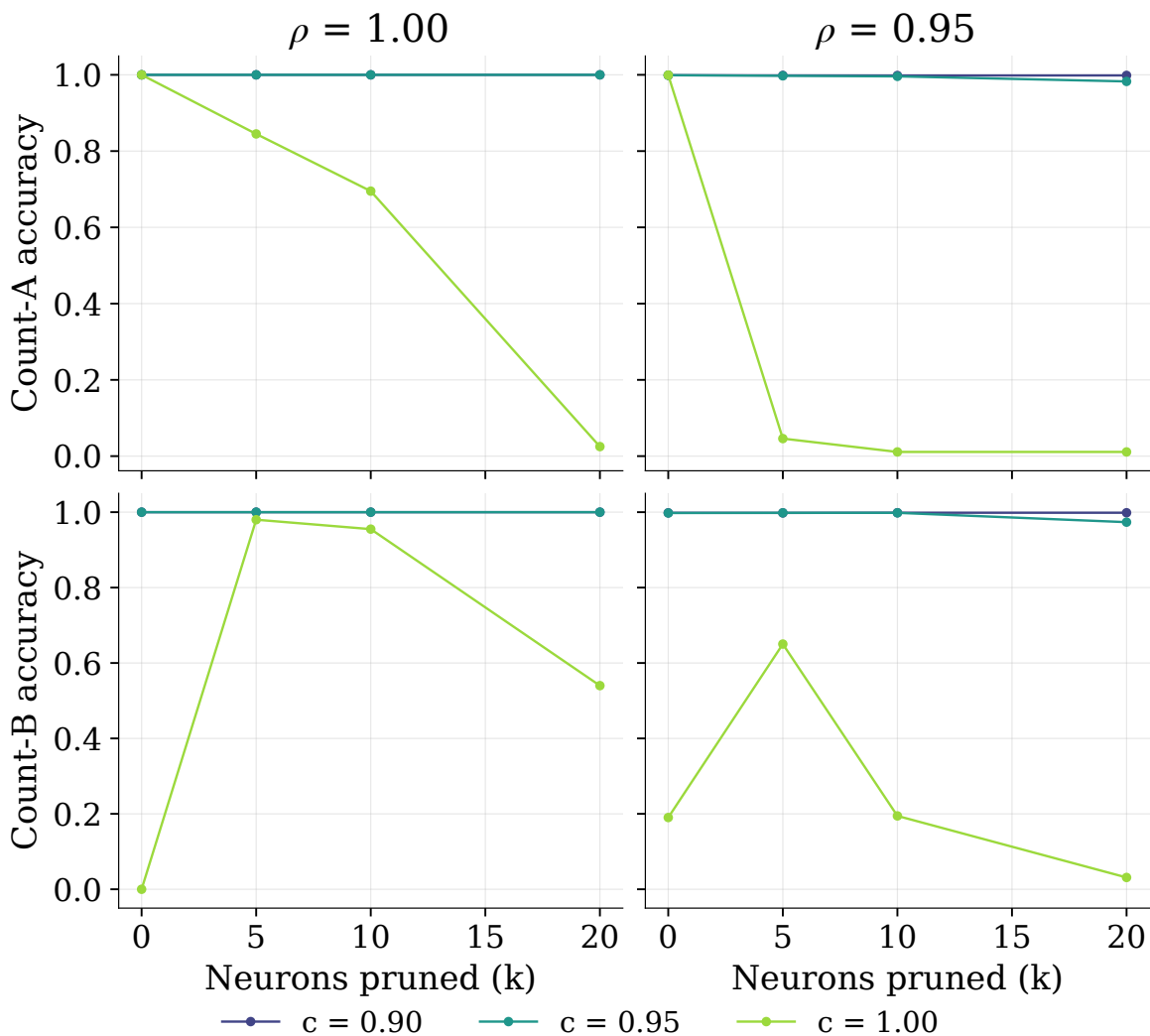


Figure 22: The resulting accuracy on Ca40 and Cb40 after pruning neurons on the PCFG setup with varying values of ρ and c . In the case of $c = 1$, pruning ~ 10 neurons recovers Cb40 accuracy while collapsing Ca40, replicating the wrapper-formation result of [14]. Reducing the concentration seems to avoid the formation of this wrapper.

Here we deviate slightly from the pruning method in the original paper, which was to choose parameters based on their gradients to most improve the performance on task Cb40;

here, we instead scan through pruning every parameter, and choose those that most improve Cb40. The reason for this choice is for lower concentrations, the model still achieves a low loss on Cb40, and using the original technique was noisy. Crucially, the new technique still reproduces the wrapper observation from the original paper.

Appendix H. Related Work

We organise related work around five themes: training order and curriculum effects (§H.1), batch diversity and data mixing (§H.2), fine-tuning behaviour and its limitations (§H.3), alignment and safety data integration (§H.4), and catastrophic forgetting and continual learning (§H.5).

H.1. Training order and curriculum effects

Curriculum learning. Bengio et al. [2] showed that presenting training examples in a meaningful order — typically from easy to hard — can improve convergence speed and generalisation. Subsequent work studies data ordering as a tool for improving *performance*; we study it as a cause of *shallowness*. The distinction is important: curriculum learning asks “what order makes the model learn best?” while we ask “what arrangement makes learning most durable?”

Out-of-context learning. Berglund et al. [3] showed that language models can learn facts stated outside their typical usage context and later deploy them at test time, and that what is retained depends on how declarative information is distributed across training. We share the observation that order matters, but isolate a specific structural variable — temporal concentration — rather than studying order in general.

Training-order recency. Krasheninnikov et al. [20] show that models linearly encode training-order recency in their activations. Our concentration sweep disentangles recency from concentration: all schedules see the target data at the same temporal position (the end of training), but at different concentrations. The fact that lower concentration produces more durable learning, despite an identical recency profile, shows that concentration — not recency — is the operative variable.

Spaced repetition. The cognitive science literature on spaced repetition [8] has long established that distributing study sessions over time produces more durable memory than massing them into a single block. Our work provides a mechanistic account of *why* this occurs in neural networks: spaced (mixed) schedules force the optimiser to find solutions that simultaneously satisfy both the target and base distributions, producing cooperative gradient directions and integrated representations.

H.2. Batch diversity and data mixing

Data mixing in pre-training. Large-scale pre-training pipelines routinely mix data from multiple domains within each batch [35]. The rationale is typically framed in terms of domain coverage and preventing catastrophic forgetting. Our work provides theoretical grounding for this practice: mixing prevents the temporal concentration that produces shallow learning. It also suggests that the *degree* of mixing matters — not just whether domains are mixed but how concentrated each domain is within the schedule.

Batch diversity. Determinantal Point Processes have been used to diversify minibatch composition [37], showing that more diverse batches can improve training. Class-balanced

sampling [9] addresses imbalanced datasets by ensuring each class is proportionally represented. These methods operate at the *within-batch* level; our concentration parameter operates at the *across-batch* level (how the target distribution is spread across training steps).

H.3. Fine-tuning behaviour and limitations

Fine-tuning shallowness and wrappers. Jain et al. [14] demonstrated that fine-tuning on PCFG tasks produces wrapper-like changes: low-rank, later-layer-localised, and easily reversible. Related work on mode connectivity [24] and skill localisation [27] similarly finds that fine-tuning tends to produce mechanistically narrow, spatially localised updates. Our work extends these findings by showing *why* wrappers form: the temporal concentration of fine-tuning data is a (and possibly *the*) causal driver. We demonstrate that the wrapper signature — opposing gradients, localised centroid shifts, low-rank weight deltas — fades monotonically as concentration decreases, establishing a continuous relationship between data scheduling and the depth of learning.

Safety fine-tuning fragility. Qi et al. [29] showed that safety fine-tuning of aligned models is easily compromised by a small number of adversarial examples, and that even benign fine-tuning on downstream tasks can degrade safety. These works demonstrate *that* safety fine-tuning is fragile; we identify a *structural cause* (concentration) and show that it is *tunable* — reducing concentration makes the same safety training more durable, without reducing its effectiveness.

LoRA and low-rank fine-tuning. LoRA [12] and related parameter-efficient methods constrain fine-tuning updates to be low-rank by design. The empirical success of LoRA provides independent structural evidence that fine-tuning naturally produces low-rank changes. Our gradient effective rank measurements confirm this: at high concentration, the gradient update is indeed low-rank. However, our results also show that this low-rank structure is a *symptom* of concentration, not an inherent requirement of the task: at lower concentration, gradient effective rank is higher and the model finds qualitatively different, more distributed solutions.

Machine unlearning. The machine unlearning literature [7] studies how to remove specific data or capabilities from trained models. Our shallowness metrics — particularly the Reversion AUC and gradient interference measures — directly inform when unlearning is expected to be easy (high concentration → shallow → easy to remove) versus difficult (low concentration → deep → harder to remove).

H.4. Alignment and safety data integration

Alignment during pre-training. Recent work has explored integrating alignment data during pre-training rather than applying it as a post-hoc fine-tuning step. Korbak et al. [18] showed that models pre-trained with human preference data embedded throughout training exhibit more robust alignment than those fine-tuned afterward; Zhou et al. [38] separately demonstrated that high-quality alignment can be achieved with surprisingly little fine-tuning data, suggesting that *where* the signal sits in the training schedule matters at least as much as *how much* of it there is. These findings are highly consistent with our hypothesis: mixing alignment data throughout training (reducing concentration) should produce deeper, more durable alignment. Our controlled experiments provide a mechanistic explanation for why this works.

Continual alignment. Several works [30] have proposed methods for maintaining alignment across multiple rounds of fine-tuning or deployment. Our results suggest that the fundamental challenge of continual alignment is a scheduling problem: if each round of alignment training is delivered at high concentration, it will produce a shallow wrapper that the next round of task training can easily undo. Interleaving alignment data with task data at each stage would produce more durable integration.

H.5. Catastrophic forgetting and continual learning

Catastrophic forgetting. The problem of catastrophic forgetting — where training on new tasks destroys performance on old tasks — [10, 25] is closely related to our work. Classical forgetting occurs when the new-task gradient overwrites the old-task parameters. Our setting inverts the direction: we study how *old-task* gradients overwrite *new-task* parameters during reversion. However, the underlying mechanism is the same — gradient interference between competing distributions. Our gradient cosine measurements directly quantify this interference.

Elastic Weight Consolidation (EWC). EWC [17] addresses catastrophic forgetting by adding a regularisation term that penalises changes to parameters important for previously learned tasks, weighted by the Fisher information matrix. EWC and related methods address forgetting by constraining the *parameter trajectory*; our work addresses it by constraining the *data schedule*. These approaches are complementary: EWC asks “which parameters should we protect?” while our work asks “how should we arrange data so that protection is unnecessary?”

Progressive neural networks and modular approaches. An alternative to overcoming forgetting is to allocate separate capacity for each task [32]. This avoids interference entirely but scales poorly. Our results suggest a middle ground: by reducing concentration (mixing tasks within batches), the model can learn to share representations across tasks without the destructive interference that concentrated training produces.

Gradient episodic memory (GEM). GEM [23] addresses continual learning by constraining gradient updates to not increase loss on previously seen tasks. This is conceptually similar to what happens naturally at low concentration in our experiments: when new-task and pre-training data are mixed, the model receives gradient signal from both distributions simultaneously, which naturally constrains the update to directions that do not increase pre-training loss. Our gradient cosine measurements (positive at low concentration, negative at high) provide direct evidence that this constraint is satisfied automatically when concentration is low.