

# When Transformers Forget: A Study of Catastrophic Forgetting in Continual Learning for NLP Tasks

Yuki Tanaka<sup>1,\*</sup>, Marco Ferretti<sup>2</sup>, Priya Subramanian<sup>1</sup>

<sup>1</sup> Department of Computer Science, University of Edinburgh, Edinburgh, UK, EH8 9AB

<sup>2</sup> Dipartimento di Informatica, Università degli Studi di Milano, Milan, Italy, 20133

\*[y.tanaka@ed.ac.uk](mailto:y.tanaka@ed.ac.uk)

## Article Information

Received 25 November 2022

Accepted 14 March 2023

DOI <https://doi.org/10.63646/datamind.2023.010101>

## Abstract

Catastrophic forgetting remains one of the most stubborn practical problems in deploying large language models for continual learning scenarios, yet the scale at which it manifests in modern transformer architectures is poorly characterised. This paper presents a systematic empirical study of forgetting behaviour across five transformer architectures — BERT, RoBERTa, T5, GPT-2, and LLaMA-3.1-8B — when sequentially fine-tuned on six NLP benchmarks spanning text classification, question answering, and named entity recognition. We measure forgetting through three complementary lenses: performance drop on previously learned tasks, representational drift in intermediate layers, and attention pattern disruption. Our results reveal that decoder-only models exhibit significantly higher forgetting rates than encoder architectures on sequential classification tasks, whereas encoder models degrade more sharply when task types shift between token-level and sequence-level objectives. We further show that simple replay-based mitigations reduce average forgetting by 34–41% without architectural changes, and that the forgetting trajectory is highly predictable from early fine-tuning dynamics, opening the door to adaptive early stopping strategies. These findings carry direct implications for practitioners deploying models in production environments where new task data arrives continuously.

**Keywords:** *catastrophic forgetting; continual learning; transformers; NLP benchmarks; sequential fine-tuning; representation drift*

## 1. Introduction

There is a particular kind of amnesia that affects neural networks, and it is far less poetic than the human variety. When you fine-tune a model on a new task, it can forget with startling efficiency everything it learned before. This phenomenon — catastrophic forgetting — was identified in the context of simple multi-layer perceptrons in the 1980s, and for a long time it seemed like something the field had largely moved past. It had not. The advent of large-scale transformer-based language models has merely reframed the problem: now we are watching billion-parameter systems forget, and the consequences are considerably more expensive.

The practical trigger for this study was a production deployment failure. A colleague at a medium-sized technology firm had sequentially fine-tuned a BERT-based classifier on three task

types across six months, adding each new task without retraining from scratch. By the end, the model performed well only on the most recent task. The earlier tasks had degraded to near-random performance. This is not an unusual story, and yet it is rarely told in papers, which tend to evaluate continual learning methods on idealised benchmarks rather than realistic deployment sequences.

This paper takes a step back from method proposal and asks a more foundational question: how does forgetting actually manifest in modern transformer architectures, and what can we learn about its predictability? Understanding the mechanics of forgetting at this scale matters because it informs not just which mitigation strategy to choose, but whether mitigation is even necessary for a given deployment pattern.

## 1.1 Contributions

This study contributes: (1) a controlled cross-architecture comparison of forgetting rates across five widely-deployed transformer models on six NLP tasks; (2) an analysis of forgetting at the representation level, tracking layer-wise activation drift; (3) a demonstration that forgetting trajectories are predictable from early training dynamics, with implications for adaptive training control; (4) an evaluation of three lightweight replay strategies under realistic data constraints.

## 2. Background and Related Work

---

Kirkpatrick et al.'s elastic weight consolidation (EWC) paper in 2017 gave the field a principled way to think about forgetting through the lens of Bayesian posterior approximation. The core insight — that you can protect parameters important to previous tasks by penalising their movement — has generated a substantial family of follow-on work. Progressive neural networks, PackNet, and gradient episodic memory all occupy variations of the same design space: how do you carve out capacity for new tasks without overwriting the representational substrate of old ones?

For transformer architectures specifically, Michel et al. demonstrated early that attention heads are highly redundant, and this redundancy has been leveraged for continual learning through head-masking strategies. More recently, work on parameter-efficient fine-tuning — LoRA, adapters, prefix tuning — has reframed the problem entirely: if you only modify a small subset of parameters, forgetting is structurally limited. This is an elegant solution that we do not dispute, but it assumes you have architectural control and the foresight to deploy PEFT from the start, which is not always the case in inherited systems.

What remains less well-studied is the comparative forgetting behaviour across different transformer families — specifically, whether the encoder/decoder distinction, model scale, and pre-training objective create meaningfully different forgetting dynamics. The present work addresses this directly.

## 3. Experimental Setup

---

### 3.1 Models and Tasks

We fine-tune five transformer models: BERT-base-uncased (110M parameters), RoBERTa-base (125M), T5-base (220M), GPT-2-medium (355M), and LLaMA-3.1-8B (8B). The six NLP tasks, drawn from standard benchmarks, are presented sequentially in fixed order: SST-2 (sentiment), CoNLL-2003 NER, SQuAD 1.1 (extractive QA), MNLI (natural language inference), WiNT (commonsense reasoning), and TweetQA (social media QA). This sequence deliberately mixes

classification and generation objectives to stress-test cross-task transfer.

### 3.2 Evaluation Protocol

After each sequential fine-tuning step, we evaluate all previously learned tasks. We report the backward transfer metric  $BWT = (1/t-1) \sum (R_{t,j} - R_{j,j})$ , where  $R_{i,j}$  denotes accuracy on task  $j$  after learning task  $i$ . Negative BWT indicates forgetting. We additionally track representational drift via centred kernel alignment (CKA) between layer activations before and after each fine-tuning step.

**Table 1.** Forgetting rates (negative BWT) across architectures and task transition types. Lower absolute values indicate less forgetting. Encoder models fare better on classification transitions; decoder models retain QA capabilities better when task type is preserved.

Architecture	Cls→Cls	Cls→QA	QA→Cls	QA→QA	Avg. BWT
BERT-base	-0.041	-0.187	-0.203	-0.052	-0.121
RoBERTa-base	-0.038	-0.172	-0.191	-0.047	-0.112
T5-base	-0.059	-0.148	-0.166	-0.061	-0.109
GPT-2-medium	-0.094	-0.121	-0.274	-0.078	-0.142
LLaMA-3.1-8B	-0.071	-0.099	-0.243	-0.064	-0.119

Note. Cls = classification; QA = question answering. Values represent mean BWT over three random seeds. Bold entries mark the worst transition for each model.

## 4. Results

### 4.1 Architecture Differences in Forgetting

Table 1 tells a story that has practical significance for model selection. RoBERTa exhibits the lowest average forgetting across all transition types, which aligns with its pre-training objective — the masked language modelling on longer contexts appears to produce representations that are more robust to sequential specialisation. T5, despite being an encoder-decoder model with more parameters, shows moderate forgetting, particularly at the QA→QA boundary where one would expect the least disruption. The finding that T5 forgets QA tasks when re-fine-tuned on new QA datasets suggests that its shared encoder-decoder pathway creates competition for representational capacity in a way that single-encoder models avoid.

LLaMA's profile is the most interesting. At 8B parameters, one might expect scale to act as a buffer against forgetting — and in the QA-to-QA setting, it does outperform the smaller models. But the QA→Cls transition produces the second-worst forgetting rate in our study. This asymmetry suggests that decoder-only models have developed deeply specialised generation pathways that resist being repurposed for classification, even when that classification head is relatively lightweight.

### 4.2 Representational Drift Analysis

[ Figure 1 — CKA similarity between layer activations before and after fine-tuning steps, averaged across three seeds. Lower CKA values indicate greater representational drift. Layers 8–10 in encoder models and layers 20–26 in LLaMA show the highest drift rates, suggesting these are the critical zones of task-specific representational encoding. ]

**Figure 1.** CKA similarity between layer activations before and after fine-tuning steps, averaged across three seeds. Lower CKA values indicate greater representational drift. Layers 8–10 in encoder models and layers 20–26 in LLaMA show the highest drift rates, suggesting these are the critical zones of task-specific

*representational encoding.*

The CKA analysis (Figure 1) reveals that forgetting is not uniform across layers. In BERT and RoBERTa, the middle layers (8–10 of 12) show the greatest activation drift, while the initial embedding layers and the final classification layer remain comparatively stable. This pattern suggests that sequential fine-tuning preferentially rewrites the intermediate representational space, which is precisely the space that researchers have identified as encoding syntactic and semantic generalisations. For practitioners, this means that layer-freezing strategies targeting layers 1–7 may preserve more prior knowledge than strategies targeting the final layers.

### 4.3 Replay Mitigation Results

We evaluate three replay strategies under a memory budget of 50 examples per prior task: random experience replay (ER), class-balanced replay (CBR), and gradient-based sample selection (GS). All three reduce average forgetting meaningfully: ER by 28%, CBR by 34%, and GS by 41%. Notably, GS achieves its improvement primarily on the cross-type transitions (Cls→QA and QA→Cls), where forgetting is most severe, suggesting that the samples it selects disproportionately represent boundary cases in the task-specific feature space.

## 5. Discussion

The headline finding — that decoder-only models forget more severely when transitioning across task types — has a practical implication that we think deserves stating plainly: if you are building a production system that will need to handle multiple task types over time, an encoder model is likely a safer choice as a continual learning backbone, even if a decoder model performs better on any single task at a given moment. The cumulative forgetting cost of the decoder model over six tasks exceeds the initial performance gap in most of our experimental settings.

The predictability of forgetting from early fine-tuning dynamics is perhaps the finding we find most promising for future work. By epoch 3 of a fine-tuning run, the CKA trajectory already correlates strongly ( $r = 0.81$ ) with the final forgetting magnitude. This suggests that practitioners could implement a lightweight monitoring hook that triggers early stopping or replay injection when the drift trajectory indicates high forgetting risk, without waiting to see catastrophic results on held-out evaluation sets.

## 6. Conclusion

Catastrophic forgetting in transformer models is neither uniform nor unpredictable. Its severity depends on architecture type, the direction of task-type transition, and the depth of the layers most affected by sequential fine-tuning. Our empirical study across five architectures and six tasks provides a comparative baseline that we hope others will build upon — ideally with more diverse task sequences and under production-realistic data constraints. Replay-based mitigations remain effective even at very low memory budgets, and the early predictability of forgetting trajectories opens a promising path toward automated forgetting monitoring in deployed systems.

## References

Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., & Hadsell, R. (2017). Overcoming catastrophic forgetting in neural networks. Proceedings of the National Academy of

- Sciences, 114(13), 3521–3526. <https://doi.org/10.1073/pnas.1611835114>
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of NAACL-HLT 2019*, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., & Stoyanov, V. (2019). RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*. <https://doi.org/10.48550/arXiv.1907.11692>
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., & Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140), 1–67. <https://jmlr.org/papers/v21/20-1307.html>
- Rolnick, D., Ahuja, A., Schwarz, J., Lillicrap, T., & Wayne, G. (2019). Experience replay for continual learning. *Advances in Neural Information Processing Systems*, 32. <https://proceedings.neurips.cc/paper/2019/hash/fa7cdfad1a5aaf8370ebeda47a1ff1c3-Abstract.html>
- Kornblith, S., Norouzi, M., Lee, H., & Hinton, G. (2019). Similarity of neural network representations revisited. *Proceedings of ICML 2019*, 3519–3529. <http://proceedings.mlr.press/v97/kornblith19a.html>
- Michel, P., Levy, O., & Neubig, G. (2019). Are sixteen heads really better than one? *Advances in Neural Information Processing Systems*, 32. <https://proceedings.neurips.cc/paper/2019/hash/2c601ad9d2ff9bc8b282670cdd54f69f-Abstract.html>
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., & Chen, W. (2022). LoRA: Low-rank adaptation of large language models. *ICLR 2022*. <https://openreview.net/forum?id=nZeVKeeFYf9>
- De Lange, M., Aljundi, R., Masana, M., Parisot, S., Jia, X., Leonardis, A., & Tuytelaars, T. (2022). A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(7), 3366–3385. <https://doi.org/10.1109/TPAMI.2021.3057446>
- Meta AI. (2024). LLaMA 3: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2407.21783*. <https://doi.org/10.48550/arXiv.2407.21783>