

---

# Continual Pre-training of MoEs: How robust is your router?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

Prior work has shown that a simple combination of replay and learning rate re-warming and re-decaying can enable the continual pre-training (CPT) of dense decoder-only transformers with minimal performance degradation compared to full re-training. In the case of decoder-only MoE transformers, however, it is unclear how the routing algorithm will impact continual pre-training performance: 1) *do the MoE transformer's routers exacerbate forgetting relative to a dense model?*; 2) *do the routers maintain a balanced load on previous distributions after CPT?*; 3) *are the same strategies applied to dense models sufficient to continually pre-train MoE LLMs?* In what follows, we conduct a large-scale ( $> 2B$  parameter switch and DeepSeek MoE LLMs trained for 600B tokens) empirical study across four MoE transformers to answer these questions. Our results establish a surprising robustness to distribution shifts for MoEs using both Sinkhorn-Balanced and Z-and-Aux-loss-balanced routing algorithms, even in MoEs continually pre-trained without replay. Moreover, we show that MoE LLMs maintain their sample efficiency (relative to a FLOP-matched dense model) during CPT and that they can match the performance of a fully re-trained MoE at a fraction of the cost.

## 1 Introduction

Sparsely-activated MoE transformers achieve significantly stronger performance than FLOP-matched dense models (e.g., dense models requiring the same amount of floating point operations (FLOPs) per forward pass). This is particularly advantageous in today's foundation model lifecycle, where a model spends the majority of its lifetime FLOPs being inferred. Many closed-source and open-source frontier language models have thus adopted an MoE architecture [12, 14, 29, 2, 16, 15]. Given the clear advantages of MoEs over dense transformers, practitioners will certainly want to update MoEs on new data as is currently done for dense transformers.

Without proper care, MoE transformers learn greedy routing strategies that overutilize certain experts, leading to poorer downstream performance and poorer accelerator utilization. During MoE pre-training, load balancing strategies are used to prevent such negative outcomes [18, 72, 8, 4, 12]. However, the load-balancing algorithms used in SOTA MoEs were not specifically designed for the non-IID data distributions encountered during continual pre-training. Adapting the router's decisions to a new distribution during CPT may compromise the balanced load on previous distributions, potentially leading to exacerbated forgetting and poor accelerator utilization. Avoiding these failure modes is critical for successfully updating MoE foundation models without the need for full re-training, but the continual pre-training of MoEs has not yet been thoroughly studied in the literature.

In this work, we fill the gap by providing a systematic study of MoE continual pre-training. Specifically, we select two popular routing algorithms and two popular MoE architectures used in state-of-the-art existing work [12, 46, 18, 8] to yield four different MoEs for our study. We then pre-train each MoE language model on 400B tokens of FineWeb and continually pre-train them on 200B tokens of code data and German web crawl data. Taking the strongest MoE architecture, we compare its

performance to full re-training baseline on both datasets. Our contributions can be summarized as follows.

- We demonstrate that a Penalty-Balanced (e.g. with Z and Aux loss) MoE following the DeepSeek architecture can successfully match the performance of a full re-training baseline, at a fraction of the cost.
- We show that MoEs using either Penalty-Balanced or Sinkhorn-balanced routing algorithms are *surprisingly* robust to distribution shifts in terms of 1) language modeling performance, 2) evaluation benchmarks, and 3) maximum routing imbalance.

## 2 Background and Related Work

This section provides a concise summary of the relevant background for our study. A more detailed version can be found in Sections A and B of the appendix.

**Continual Pre-training of LLMs.** Continual pre-training (CPT) extends pre-training to one or more new distributions. Concretely, continual pre-training occurs when a model is trained on a sequence of datasets  $\mathcal{D}_0, \mathcal{D}_1, \dots, \mathcal{D}_N$  with different distributions,  $N \geq 2$ , and each dataset is sufficiently large (e.g.,  $> 100\text{B}$  tokens in the case of language). Recently, Ibrahim et al. [26] established that CPT LLMs can match the performance of full re-training by simply repeating the pre-training schedule (cosine annealing) for CPT and replaying previous data. However, if one has control over the initial pre-training, it is possible to further improve CPT by using an infinite learning rate schedule [28].

**Sparingly-activated MoE transformers** differ from their dense counterparts by dynamically routing tokens in a sequence. Algorithms for dynamically selecting among experts, known as routing algorithms, are therefore central to MoEs and may play a crucial role during distribution shifts. In this work, we focus on studying two prominent Top- $k$  routing algorithms from recent state-of-the-art (SOTA) works: **Penalty-Balanced Top- $k$**  (PBT $k$ ) routing [58, 12, 72, 18] and **Sinkhorn-Balanced Top- $k$**  (SBT $k$ ) routing [8, 4].

## 3 Empirical Study

To study large-scale continual pre-training of MoE LLMs, we construct two settings using three datasets: FineWeb [49], the Stack [34], and German Common Crawl [1]. We initially pre-train all models on FineWeb for 400B tokens (task 1) to mimic open and closed source models often pre-trained on large-scale web-scraped English data. Subsequently, we continually pre-train these base models on 200B tokens of Code or German data (task 2) using infinite learning rate schedules and replay (30% & 40%, respectively) to mitigate forgetting. Using these settings, our study compares the performance of the dense and MoE architectures outlined below. See section A.3 for extended training details.

**FLOP-matched Dense Baseline.** We select a 24 layer 570M parameter dense decoder-only transformer following the Llama3 architecture (except we use GeLU activations) and using the Llama3 tokenizer [17] (see Sec. E for details).

**Granular MoEs.** Given the recent popularity and strong performance of DeepSeek MoEs [14, 12, 16, 15], we include an MoE architecture that activates multiple granular experts and a shared expert. Specifically, each granular MoE has  $E = 31$  total routed experts,  $K = 3$  active experts, and 1 shared expert. This model follows the same Llama3 architecture as the dense model described above. Notably, its experts are GEGLU FFNs with an intermediate size that is  $\frac{1}{4}$  the size used in the dense model. We train two granular MoEs utilizing the Penalty-Balanced and Sinkhorn-Balanced Top- $k$  routing algorithms, respectively. We do not drop tokens.

**Switch MoEs.** Given the historical use of full-sized FFNs in MoEs, our study also includes an architecture similar to [29, 18] with full-sized experts and no shared expert. We refer to these as Switch MoEs and also train two utilizing the Penalty-Balanced and Sinkhorn-Balanced routing algorithms, respectively. Each switch MoE has  $E = 8$  total routed experts,  $K = 1$  active experts, and no shared expert. This model follows the same Llama3 architecture as the dense model described above. Notably, its experts are GEGLU FFNs of the same size as the dense model’s FFNs. We do not drop tokens.

## 4 Results

**Validation Loss.** Table 1 reports validation loss (e.g., log perplexity on a fixed held-out validation set) results for the main models in our study, while extended results are reported in Table 5 of the appendix.

Table 1: **Aggregated benchmark results.** MoEs consistently outperform FLOP-matched dense baselines and exhibit less forgetting w.r.t. validation loss. Compared to the re-training baseline (blue), MoEs and the dense model match or exceed their performance. These results show MoEs adapt as well as dense models but forget less, likely due to their larger parameter count. All validation losses report log perplexity on a held-out validation set, forgetting (equivalent to backward transfer in [41]) is calculated using validation loss, and downstream tasks report accuracy.

Training Tokens	Model	Final Validation Loss (↓)					Downstream Evals. (↑)		
		FineWeb	Stack	German	Forgetting	AVG	English	German	Stack
400B FineWeb	Dense Baseline	2.881	4.028	3.741	–	–	49.84%	23.54%	0.00%
	SB Switch MoE	2.711	3.861	3.495	–	–	54.14%	23.11%	0.00%
	PB Switch MoE	2.699	3.872	3.451	–	–	54.45%	23.37%	0.00%
	SB Granular MoE	2.664	3.690	3.404	–	–	<b>55.71%</b>	22.83%	0.00%
	PB Granular MoE	2.653	3.715	3.370	–	–	55.59%	23.40%	0.00%
400B FineWeb → 200B Stack (30% Replay)	Dense Baseline	2.939	1.026	–	0.059	1.982	49.21%	–	7.37%
	SB Switch MoE	2.757	0.944	–	0.046	1.850	51.76%	–	<b>9.09%</b>
	PB Switch MoE	2.749	0.945	–	0.050	1.847	52.59%	–	8.22%
	SB Granular MoE	2.708	<b>0.925</b>	–	<b>0.044</b>	1.816	53.51%	–	7.45%
	PB Granular MoE	2.699	<b>0.924</b>	–	0.046	1.811	53.70%	–	7.81%
400B FineWeb ∪ 200B Stack	Dense Baseline	2.866	1.050	–	–	1.958	49.57%	–	3.76%
	PB Granular MoE	<b>2.630</b>	0.935	–	–	<b>1.782</b>	54.79%	–	7.44%
400B FineWeb → 200B German (40% Replay)	Dense Baseline	2.946	–	1.367	0.066	2.157	48.15%	25.27%	–
	SB Switch MoE	2.749	–	1.142	0.039	1.946	51.99%	27.57%	–
	PB Switch MoE	2.741	–	1.129	0.042	1.935	51.25%	26.50%	–
	SB Granular MoE	2.701	–	1.118	<b>0.037</b>	1.910	53.35%	<b>28.57%</b>	–
	PB Granular MoE	2.690	–	<b>1.099</b>	<b>0.037</b>	<b>1.895</b>	53.61%	27.65%	–
400B FineWeb ∪ 200B German	Dense Baseline	2.938	–	1.390	–	2.164	48.42%	25.45%	–
	PB Granular MoE	2.669	–	1.120	–	<b>1.895</b>	53.94%	27.59%	–

We observe that all MoEs outperform the FLOP-matched dense baseline during pre-training and CPT. Within the MoEs, we observe that PBT $k$  MoEs consistently outperform SBT $k$  MoEs and that Granular MoEs consistently outperform switch MoEs across pre-training and CPT. These findings are consistent with the literature on Granular MoEs [43, 12], but we believe this is the first time that SBT $k$  routing has been shown to underperform PBT $k$  routing. Since the PBT $k$  Granular MoE achieves the best performance, we use it as our full re-training baseline. Compared to full re-training, our Granular CPT MoEs consistently have higher FineWeb validation loss but achieve lower downstream validation loss with a similar average validation loss. Similar results are found when comparing the CPT dense baseline to its full re-training counterpart. These results demonstrate that the continual learning abilities of MoEs w.r.t. validation loss are on par with dense models for adaptation and are slightly superior in terms of forgetting, likely due to their larger total parameter count.

**English Evaluation results.** Table 1 presents average accuracy, while Table 3 details per-benchmark results. We select benchmarks where models of our scale (570M active parameters) achieve non-trivial accuracy to maximize signal. Each model is evaluated zero-shot on benchmarks covering Commonsense Reasoning, Reading Comprehension, Scientific Question Answering, and Math: HellaSwag [67], Winogrande [56], PIQA [5], ARC-Easy, ARC-Challenge [9], SWAG [66], LAMBADA (OpenAI) [60], SciQ [31], PubMedQA [30], and MathQA [3] (see Sec. C.4 for details). We find that models trained solely on FineWeb outperform all others, including full re-training baselines. Granular MoEs surpass switch MoEs. CPT models trained on Stack perform similarly to those trained on German. Compared to full re-training, CPT models achieve nearly identical results (within  $\sim 1\%$ ). The dense baseline also matches its full re-training counterpart, indicating that MoEs have similar continual learning abilities on pre-training evaluations while benefiting from improved sample efficiency.

**German Evaluation results.** Table 1 shows average German evaluation performance, while table 4 of the appendix provides a per-benchmark breakdown. We use GPT-3–translated German versions of HellaSwag, ARC-Challenge, and TruthfulQA, evaluating each zero-shot [50]. German-trained models outperform English-only ones, and German-trained MoEs surpass the FLOP-matched dense baseline. Among MoEs, modules using the same training tokens perform similarly. CPT MoEs and the full re-training baseline differ by  $< 1\%$  accuracy, with no clear winner. The dense baseline also performs comparably to full re-training, demonstrating that the continual learning abilities of MoEs w.r.t. German evaluations are on par with dense models while benefiting from improved sample efficiency.

**Code Evaluation results.** Table 1 presents average Code evaluation performance, while Table 2 of the appendix provides a pass@ $k$  breakdown ( $k \in \{1, 10, 50, 100, 150, 200\}$ ). Our models are evaluated on Python code-generation tasks from HumanEval [7], as Python is well-represented

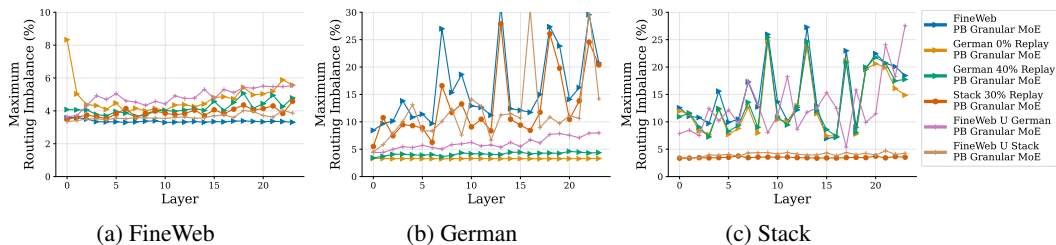


Figure 1: **Layer-wise Maximum Routing Imbalance (MRI) for Granular MoEs.** We report MRI (eq. 1) on each dataset’s 20M-token test set. MRI is consistently lower for Penalty-Balanced MoEs than Sinkhorn-Balanced MoEs. Continual pre-training on FineWeb incurs minimal MRI increase, even with 0% replay. MoEs are most unbalanced with out-of-distribution data (e.g., non-German models in (b) and non-code models in (c)).

in our Stack CPT dataset (Table 9). English-trained models can not solve any problem, whereas stack-trained models achieve non-trivial accuracy. Unlike for other performance metrics, CPT switch MoEs slightly outperform their granular counterparts. Compared to full re-training, all CPT MoEs perform marginally better, while the CPT dense model exceeds its baseline by over 3%. We attribute this unexpected improvement to evaluation noise and training variance, given the models’ similar validation loss. These results suggest MoEs match dense models in continual learning for code evaluation when accounting for MoEs’ improved sample efficiency.

**Routing imbalance during and after continual pre-training.** Figure 1 shows the layer-wise Maximum Routing Imbalance (MRI) for Granular MoEs across FineWeb (a), German (b), and stack (c), while Figure 16 reports MRI for all MoEs. We include a 0% replay baseline for each MoE CPT on German to highlight replay’s impact on MRI. In subfigure (a), Penalty-Balanced MoEs consistently have lower MRI than Sinkhorn-Balanced MoEs across all architectures, and granular MoEs exhibit lower and more stable MRI within architectures. On FineWeb, continual pre-training causes only a slight MRI increase relative to the pre-trained checkpoint, even for the 0% replay model, except in its first layer. Interestingly, all German-trained MoEs show higher MRI on FineWeb than their Stack-trained counterparts, with the full re-training baseline, surprisingly, having the highest. This suggests there may be more routing interference between English and German datasets and that continual pre-training may help reduce MRI across distributions, possibly due to the use of *CosineInf* vs. *Cosine Annealing* schedules for CPT and re-training, respectively. Granular MoEs also reduce routing imbalance on German (b) and Stack (c) datasets. MoEs become most unbalanced with out-of-distribution data (e.g., non-German models in (b) and non-code models in (c)). Similar trends hold for Switch MoEs, with an additional finding: high MRI is common in early layers of Switch MoEs, independent of training/testing distributions, unlike Granular MoEs. These results show that PBT $k$  and SBT $k$  MoEs are robust to distribution shifts w.r.t. MRI and can even outperform re-training baselines, suggesting that continually pre-training MoEs should have no negative impact on inference latency.

*In summary, we find that, across three measures of performance, MoEs continually pre-trained with replay and infinite LR schedules can match the performance of a full re-training baseline and, thus, they have similar CPT abilities to a FLOP-matched dense baseline without any inhibition from their routers. Moreover, we show that continually pre-training MoEs has no negative impact on MRI compared to re-training.*

## 5 Conclusion

We have conducted a comprehensive empirical study on the continual pre-training of decoder-only MoE transformer language models. Our large-scale experiments, involving 2B parameter MoEs trained on 600B tokens, demonstrate that both Penalty-Balanced (PBT $k$ ) and Sinkhorn-Balanced (SBT $k$ ) routing algorithms exhibit surprising system-level resilience to distribution shifts, maintaining balanced loads as measured by the novel Maximum Routing Imbalance metric. We established that MoEs preserve their sample efficiency advantage over FLOP-matched dense models during CPT and that, when using infinite LR schedules and replay, a Granular PBT $k$  MoEs can match the performance of fully re-trained baselines across German and Code transitions, at a fraction of the computational cost.

## References

- [1] Julien Abadji, Pedro Javier Ortiz Suárez, Laurent Romary, and Benoît Sagot. Towards a cleaner document-oriented multilingual crawled corpus. In Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, Hélène Mazo, Jan Odiijk, and Stelios Piperidis, editors, *Proceedings of the Thirteenth Language Resources and Evaluation Conference, LREC 2022, Marseille, France, 20-25 June 2022*, pages 4344–4355. European Language Resources Association, 2022. URL <https://aclanthology.org/2022.lrec-1.463>.
- [2] Marah I Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat S. Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Parul Chopra, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Dan Iter, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Chen Liang, Weishung Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Xia Song, Masahiro Tanaka, Xin Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Michael Wyatt, Can Xu, Jiahang Xu, Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone. *CoRR*, abs/2404.14219, 2024. doi: 10.48550/ARXIV.2404.14219. URL <https://doi.org/10.48550/arXiv.2404.14219>.
- [3] Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. Mathqa: Towards interpretable math word problem solving with operation-based formalisms, 2019.
- [4] Quentin Anthony, Yury Tokpanov, Paolo Glorioso, and Beren Millidge. Blackmamba: Mixture of experts for state-space models. *CoRR*, abs/2402.01771, 2024. URL <https://doi.org/10.48550/arXiv.2402.01771>.
- [5] Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language, 2019.
- [6] Weilin Cai, Juyong Jiang, Fan Wang, Jing Tang, Sunghun Kim, and Jiayi Huang. A survey on mixture of experts. *CoRR*, abs/2407.06204, 2024. URL <https://doi.org/10.48550/arXiv.2407.06204>.
- [7] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgén Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- [8] Aidan Clark, Diego de Las Casas, Aurelia Guy, Arthur Mensch, Michela Paganini, Jordan Hoffmann, Bogdan Damoc, Blake A. Hechtman, Trevor Cai, Sebastian Borgeaud, George van den Driessche, Eliza Rutherford, Tom Hennigan, Matthew J. Johnson, Albin Cassirer, Chris Jones, Elena Buchatskaya, David Budden, Laurent Sifre, Simon Osindero, Oriol Vinyals, Marc’Aurelio Ranzato, Jack W. Rae, Erich Elsen, Koray Kavukcuoglu, and Karen Simonyan. Unified scaling laws for routed language models. In Kamalika Chaudhuri, Stefanie Jegelka,



- 222 Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International Conference*  
 223 *on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume  
 224 162 of *Proceedings of Machine Learning Research*, pages 4057–4086. PMLR, 2022. URL  
 225 <https://proceedings.mlr.press/v162/clark22a.html>.
- 226 [9] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick,  
 227 and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning  
 228 challenge, 2018.
- 229 [10] Ronan Collobert, Yoshua Bengio, and Samy Bengio. Scaling large learning problems with  
 230 hard parallel mixtures. *Int. J. Pattern Recognit. Artif. Intell.*, 17(3):349–365, 2003. URL  
 231 <https://doi.org/10.1142/S0218001403002411>.
- 232 [11] Andrea Cossu, Tinne Tuytelaars, Antonio Carta, Lucia Passaro, Vincenzo Lomonaco, and  
 233 Davide Bacciu. Continual pre-training mitigates forgetting in language and vision, 2022. URL  
 234 <https://arxiv.org/abs/2205.09357>.
- 235 [12] Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu, Huazuo Gao, Deli Chen, Jiashi  
 236 Li, Wangding Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli Luo,  
 237 Chong Ruan, Zhifang Sui, and Wenfeng Liang. Deepseekmoe: Towards ultimate expert  
 238 specialization in mixture-of-experts language models. In Lun-Wei Ku, Andre Martins, and  
 239 Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for*  
 240 *Computational Linguistics (Volume 1: Long Papers)*, *ACL 2024, Bangkok, Thailand, August*  
 241 *11-16, 2024*, pages 1280–1297. Association for Computational Linguistics, 2024. URL [https://](https://aclanthology.org/2024.acl-long.70)  
 242 [aclanthology.org/2024.acl-long.70](https://aclanthology.org/2024.acl-long.70).
- 243 [13] MohammadReza Davari, Nader Asadi, Sudhir Mudur, Rahaf Aljundi, and Eugene Belilovsky.  
 244 Probing representation forgetting in supervised and unsupervised continual learning. In *Pro-*  
 245 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
 246 16712–16721, 2022.
- 247 [14] DeepSeek-AI, Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu,  
 248 Y. Wu, Yukun Li, Huazuo Gao, Shirong Ma, Wangding Zeng, Xiao Bi, Zihui Gu, Hanwei  
 249 Xu, Damai Dai, Kai Dong, Liyue Zhang, Yishi Piao, Zhibin Gou, Zhenda Xie, Zhewen Hao,  
 250 Bingxuan Wang, Junxiao Song, Deli Chen, Xin Xie, Kang Guan, Yuxiang You, Aixin Liu,  
 251 Qiushi Du, Wenjun Gao, Xuan Lu, Qinyu Chen, Yaohui Wang, Chengqi Deng, Jiashi Li,  
 252 Chenggang Zhao, Chong Ruan, Fuli Luo, and Wenfeng Liang. Deepseek-coder-v2: Breaking  
 253 the barrier of closed-source models in code intelligence. *CoRR*, abs/2406.11931, 2024. URL  
 254 <https://doi.org/10.48550/arXiv.2406.11931>.
- 255 [15] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin  
 256 Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu,  
 257 Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan  
 258 Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang,  
 259 Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli  
 260 Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng  
 261 Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li,  
 262 Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian  
 263 Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean  
 264 Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan  
 265 Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian,  
 266 Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong  
 267 Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan  
 268 Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan,  
 269 S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang,  
 270 Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao  
 271 Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin  
 272 Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin,  
 273 Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun,  
 274 Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong  
 275 Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yudian Wang, Yue

- 276 Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang  
277 Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian  
278 Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu,  
279 Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu,  
280 Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu,  
281 Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via  
282 reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- 283 [16] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu,  
284 Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian  
285 Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao,  
286 Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang,  
287 Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo,  
288 Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong  
289 Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean  
290 Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li,  
291 Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian,  
292 Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du,  
293 R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu  
294 Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu,  
295 Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng  
296 Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng,  
297 Wanbiao Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang,  
298 X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen,  
299 Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang,  
300 Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi  
301 Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei,  
302 Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng  
303 Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying  
304 He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo,  
305 Yu Wu, Yuan Ou, Yuchen Zhu, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha,  
306 Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou,  
307 Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang,  
308 Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong  
309 Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu,  
310 Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report,  
311 2025. URL <https://arxiv.org/abs/2412.19437>.
- 312 [17] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle,  
313 Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony  
314 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,  
315 Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière,  
316 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi,  
317 Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne  
318 Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle  
319 Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano,  
320 Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily  
321 Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee,  
322 Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell,  
323 Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra,  
324 Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana  
325 Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny  
326 Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu,  
327 Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng  
328 Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield,  
329 Kevin Stone, and et al. The llama 3 herd of models. *CoRR*, abs/2407.21783, 2024. URL  
330 <https://doi.org/10.48550/arXiv.2407.21783>.
- 331 [18] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion  
332 parameter models with simple and efficient sparsity. *J. Mach. Learn. Res.*, 23:120:1–120:39,

2022. URL <http://jmlr.org/papers/v23/21-0998.html>.
- [19] Robert M. French. Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, 3(4):128–135, 1999. ISSN 13646613. doi: 10.1016/S1364-6613(99)01294-2. URL <https://www.sciencedirect.com/science/article/abs/pii/S1364661399012942>.
- [20] Trevor Gale, Deepak Narayanan, Cliff Young, and Matei Zaharia. MegaBlocks: Efficient Sparse Training with Mixture-of-Experts. *Proceedings of Machine Learning and Systems*, 5, 2023.
- [21] Saurabh Garg, Mehrdad Farajtabar, Hadi Pouransari, Raviteja Vemulapalli, Sachin Mehta, Oncel Tuzel, Vaishaal Shankar, and Fartash Faghri. Tic-clip: Continual training of clip models. *arXiv preprint arXiv:2310.16226*, 2023. URL <https://arxiv.org/abs/2310.16226>.
- [22] Nikolas Gritsch, Qizhen Zhang, Acyr Locatelli, Sara Hooker, and Ahmet Üstün. Nexus: Specialization meets adaptability for efficiently training mixture of experts, 2024. URL <https://arxiv.org/abs/2408.15901>.
- [23] Kshitij Gupta, Benjamin Thérien, Adam Ibrahim, Mats Leon Richter, Quentin Gregory Anthony, Eugene Belilovsky, Irina Rish, and Timothée Lesort. Continual pre-training of large language models: How to re-warm your model? In *Workshop on Efficient Systems for Foundation Models @ ICML2023*, 2023. URL <https://openreview.net/forum?id=pg7PUJe0Tl>.
- [24] Xu Owen He. Mixture of A million experts. *CoRR*, abs/2407.04153, 2024. URL <https://doi.org/10.48550/arXiv.2407.04153>.
- [25] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models. *CoRR*, abs/2203.15556, 2022. URL <https://doi.org/10.48550/arXiv.2203.15556>.
- [26] Adam Ibrahim, Benjamin Thérien, Kshitij Gupta, Mats Leon Richter, Quentin Gregory Anthony, Eugene Belilovsky, Timothée Lesort, and Irina Rish. Simple and scalable strategies to continually pre-train large language models. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=DimPeeCxK0>.
- [27] Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. Adaptive mixtures of local experts. *Neural Comput.*, 3(1):79–87, 1991. URL <https://doi.org/10.1162/neco.1991.3.1.79>.
- [28] Paul Janson, Vaibhav Singh, Paria Mehrbod, Adam Ibrahim, Irina Rish, Eugene Belilovsky, and Benjamin Thérien. Beyond cosine decay: On the effectiveness of infinite learning rate schedule for continual pre-training. *arXiv preprint arXiv:2503.02844*, 2025.
- [29] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Léo Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts. *CoRR*, abs/2401.04088, 2024. URL <https://doi.org/10.48550/arXiv.2401.04088>.
- [30] Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. Pubmedqa: A dataset for biomedical research question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2567–2577, 2019.
- [31] Matt Gardner Johannes Welbl, Nelson F. Liu. Crowdsourcing multiple choice science questions. *arXiv:1707.06209v1*, 2017.



- [32] Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In Regina Barzilay and Min-Yen Kan, editors, *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1147. URL <https://aclanthology.org/P17-1147>.
- [33] Paul Knopp and Richard Sinkhorn. Concerning nonnegative matrices and doubly stochastic matrices. *Pacific Journal of Mathematics*, 21(2):343 – 348, 1967.
- [34] Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Yacine Jernite, Margaret Mitchell, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Dzmitry Bahdanau, Leandro von Werra, and Harm de Vries. The stack: 3 TB of permissively licensed source code. *Trans. Mach. Learn. Res.*, 2023, 2023. URL <https://openreview.net/forum?id=pxpbTdUEpD>.
- [35] Aran Komatsuzaki, Joan Puigcerver, James Lee-Thorp, Carlos Riquelme Ruiz, Basil Mustafa, Joshua Ainslie, Yi Tay, Mostafa Dehghani, and Neil Houlsby. Sparse upcycling: Training mixture-of-experts from dense checkpoints. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/forum?id=T5nUQDrM4u>.
- [36] Dmitry Lepikhin, Hyoungho Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=qrwe7XHTmYb>.
- [37] Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. BASE layers: Simplifying training of large, sparse models. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 6265–6274. PMLR, 2021. URL <http://proceedings.mlr.press/v139/lewis21a.html>.
- [38] Liyuan Liu, Jianfeng Gao, and Weizhu Chen. Sparse backpropagation for moe training. *CoRR*, abs/2310.00811, 2023. URL <https://doi.org/10.48550/arXiv.2310.00811>.
- [39] Liyuan Liu, Young Jin Kim, Shuohang Wang, Chen Liang, Yelong Shen, Hao Cheng, Xiaodong Liu, Masahiro Tanaka, Xiaoxia Wu, Wenxiang Hu, Vishrav Chaudhary, Zeqi Lin, Chengruidong Zhang, Jilong Xue, Hany Awadalla, Jianfeng Gao, and Weizhu Chen. GRIN: gradient-informed moe. *CoRR*, abs/2409.12136, 2024. URL <https://doi.org/10.48550/arXiv.2409.12136>.
- [40] Zeyu Liu, Tim Dettmers, Xi Lin, Veselin Stoyanov, and Xian Li. Towards A unified view of sparse feed-forward network in pretraining large language model. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 15038–15061. Association for Computational Linguistics, 2023. URL <https://doi.org/10.18653/v1/2023.emnlp-main.930>.
- [41] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 6467–6476, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/f87522788a2be2d171666752f97ddeb-Abstract.html>.
- [42] Ilya Loshchilov and Frank Hutter. SGDR: stochastic gradient descent with warm restarts. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=Skq89Scxx>.

- [43] Jan Ludziejewski, Jakub Krajewski, Kamil Adamczewski, Maciej Pióro, Michal Krutul, Szymon Antoniak, Kamil Ciebiera, Krystian Król, Tomasz Odrzygóźdź, Piotr Sankowski, Marek Cygan, and Sebastian Jaszczur. Scaling laws for fine-grained mixture of experts. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. URL <https://openreview.net/forum?id=yoqdl9nCRs>.
- [44] Sanket Vaibhav Mehta, Darshan Patil, Sarath Chandar, and Emma Strubell. An empirical investigation of the role of pre-training in lifelong learning. *J. Mach. Learn. Res.*, 24:214:1–214:50, 2023. URL <http://jmlr.org/papers/v24/22-0496.html>.
- [45] Seyed-Iman Mirzadeh, Arslan Chaudhry, Dong Yin, Huiyi Hu, Razvan Pascanu, Dilan Görür, and Mehrdad Farajtabar. Wide neural networks forget less catastrophically. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 15699–15717. PMLR, 2022. URL <https://proceedings.mlr.press/v162/mirzadeh22a.html>.
- [46] Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, Yuling Gu, Shane Arora, Akshita Bhagia, Dustin Schwenk, David Wadden, Alexander Wettig, Binyuan Hui, Tim Dettmers, Douwe Kiela, Ali Farhadi, Noah A. Smith, Pang Wei Koh, Amanpreet Singh, and Hannaneh Hajishirzi. Olmoe: Open mixture-of-experts language models. *CoRR*, abs/2409.02060, 2024. URL <https://doi.org/10.48550/arXiv.2409.02060>.
- [47] Ashwinee Panda, Vatsal Baherwani, Zain Sarwar, Benjamin Thérien, Stephen Rawls, Sambit Sahu, Supriyo Chakraborty, and Tom Goldstein. Dense backpropagation improves routing for sparsely-gated mixture-of-experts. In *Workshop on Machine Learning and Compression, NeurIPS 2024*, 2024. URL <https://openreview.net/forum?id=9g285TLM8>.
- [48] Jupinder Parmar, Sanjeev Satheesh, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. Reuse, don’t retrain: A recipe for continued pretraining of language models. *CoRR*, abs/2407.07263, 2024. URL <https://doi.org/10.48550/arXiv.2407.07263>.
- [49] Guilherme Penedo, Hynek Kydlicek, Loubna Ben Allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro von Werra, and Thomas Wolf. The fineweb datasets: Decanting the web for the finest text data at scale. *CoRR*, abs/2406.17557, 2024. URL <https://doi.org/10.48550/arXiv.2406.17557>.
- [50] Björn Plüster. German Benchmark Datasets, 8 2023. URL <https://github.com/bjoernpl/GermanBenchmark>.
- [51] Jack W. Rae, Katie Millican, Siddhant M. Jayakumar, David Menick, Asya Berglund, Tom Hennigan, Roman Ring, Mandy Korpusik, Matthew Hechtman, Jacob Hilton, John S. Garcia, James Norman, Sasha Borgeaud, Trevor Cai, Jordan Hoffmann, Katarzyna Krawczyk, Arthur Mensch, Thomas Scialom, Eric Alford, Jordan D. L. Ho, Daniel Hesslow, Thomas Gunter, Jason Phang, Beren Millidge, Fan Yang, Marie-Anne Lachaux, Lorrayne de Souza Schmerling, Nat McAleese, Heidy Khlaaf, Simon Osindero, Oriol Vinyals, Karol Hausman, Laurent Sifre, Andrew M. Dai, Geoffrey Irving, Michael C. Mozer, Jeff Dean, and Koray Kavukcuoglu. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021. URL <https://arxiv.org/abs/2112.11446>.
- [52] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: memory optimizations toward training trillion parameter models. In Christine Cuicchi, Irene Qualters, and William T. Kramer, editors, *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC 2020, Virtual Event / Atlanta, Georgia, USA, November 9-19, 2020*, page 20. IEEE/ACM, 2020. URL <https://doi.org/10.1109/SC41405.2020.00024>.
- [53] Samyam Rajbhandari, Conglong Li, Zhewei Yao, Minjia Zhang, Reza Yazdani Aminabadi, Ammar Ahmad Awan, Jeff Rasley, and Yuxiong He. DeepSpeed-moe: Advancing mixture-of-experts inference and training to power next-generation AI scale. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International*

- 483 *Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA,*  
 484 *volume 162 of Proceedings of Machine Learning Research*, pages 18332–18346. PMLR, 2022.  
 485 URL <https://proceedings.mlr.press/v162/rajbhandari22a.html>.
- 486 [54] Vinay Venkatesh Ramasesh, Aitor Lewkowycz, and Ethan Dyer. Effect of scale on catas-  
 487 trophic forgetting in neural networks. In *The Tenth International Conference on Learning*  
 488 *Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. URL  
 489 [https://openreview.net/forum?id=GhVS8\\_yPeEa](https://openreview.net/forum?id=GhVS8_yPeEa).
- 490 [55] Stephen Roller, Sainbayar Sukhbaatar, Arthur Szlam, and Jason Weston. Hash layers for large  
 491 sparse models. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang,  
 492 and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing Systems 34:*  
 493 *Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December*  
 494 *6-14, 2021, virtual*, pages 17555–17566, 2021. URL [https://proceedings.neurips.cc/](https://proceedings.neurips.cc/paper/2021/hash/92bf5e6240737e0326ea59846a83e076-Abstract.html)  
 495 [paper/2021/hash/92bf5e6240737e0326ea59846a83e076-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/92bf5e6240737e0326ea59846a83e076-Abstract.html).
- 496 [56] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An  
 497 adversarial winograd schema challenge at scale, 2019.
- 498 [57] Thomas Scialom, Tuhin Chakrabarty, and Smaranda Muresan. Fine-tuned language models are  
 499 continual learners. In *Proceedings of the 2022 Conference on Empirical Methods in Natural*  
 500 *Language Processing*, pages 6107–6122, 2022.
- 501 [58] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc V. Le, Geoffrey E.  
 502 Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-  
 503 of-experts layer. In *5th International Conference on Learning Representations, ICLR 2017,*  
 504 *Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL  
 505 <https://openreview.net/forum?id=B1ckMDqlg>.
- 506 [59] Yikang Shen, Zhen Guo, Tianle Cai, and Zengyi Qin. Jetmoe: Reaching llama2 performance  
 507 with 0.1m dollars. *CoRR*, abs/2404.07413, 2024. URL [https://doi.org/10.48550/arXiv.](https://doi.org/10.48550/arXiv.2404.07413)  
 508 [2404.07413](https://doi.org/10.48550/arXiv.2404.07413).
- 509 [60] Shane Storks, Qiaozi Gao, and Joyce Yue Chai. Recent advances in natural language inference:  
 510 A survey of benchmarks, resources, and approaches. *arXiv: Computation and Language*, 2019.  
 511 URL <https://api.semanticscholar.org/CorpusID:213613608>.
- 512 [61] Sainbayar Sukhbaatar, Olga Golovneva, Vasu Sharma, Hu Xu, Xi Victoria Lin, Baptiste Rozière,  
 513 Jacob Kahn, Daniel Li, Wen-tau Yih, Jason Weston, and Xian Li. Branch-train-mix: Mixing  
 514 expert llms into a mixture-of-experts LLM. *CoRR*, abs/2403.07816, 2024. URL [https:](https://doi.org/10.48550/arXiv.2403.07816)  
 515 [//doi.org/10.48550/arXiv.2403.07816](https://doi.org/10.48550/arXiv.2403.07816).
- 516 [62] Qwen Team. Qwen1.5-moe: Matching 7b model performance with 1/3 activated parameters",  
 517 February 2024. URL <https://qwenlm.github.io/blog/qwen-moe/>.
- 518 [63] Lean Wang, Huazuo Gao, Chenggang Zhao, Xu Sun, and Damai Dai. Auxiliary-loss-free  
 519 load balancing strategy for mixture-of-experts. *CoRR*, abs/2408.15664, 2024. URL [https:](https://doi.org/10.48550/arXiv.2408.15664)  
 520 [//doi.org/10.48550/arXiv.2408.15664](https://doi.org/10.48550/arXiv.2408.15664).
- 521 [64] Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual  
 522 learning: Theory, method and application. *IEEE Trans. Pattern Anal. Mach. Intell.*, 46(8):  
 523 5362–5383, 2024. URL <https://doi.org/10.1109/TPAMI.2024.3367329>.
- 524 [65] Zihan Wang, Deli Chen, Damai Dai, Runxin Xu, Zhuoshu Li, and Y. Wu. Let the expert stick to  
 525 his last: Expert-specialized fine-tuning for sparse architectural large language models. *CoRR*,  
 526 abs/2407.01906, 2024. URL <https://doi.org/10.48550/arXiv.2407.01906>.
- 527 [66] Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. SWAG: A large-scale adversarial  
 528 dataset for grounded commonsense inference. In Ellen Riloff, David Chiang, Julia Hockenmaier,  
 529 and Jun’ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in*  
 530 *Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages  
 531 93–104. Association for Computational Linguistics, 2018.

- 532 [67] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a  
533 machine really finish your sentence?, 2019.
- 534 [68] Xiaofeng Zhang, Yikang Shen, Zeyu Huang, Jie Zhou, Wenge Rong, and Zhang Xiong. Mixture  
535 of attention heads: Selecting attention heads per token. In Yoav Goldberg, Zornitsa Kozareva,  
536 and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural  
537 Language Processing*, Abu Dhabi, United Arab Emirates, December 2022. Association for  
538 Computational Linguistics. URL <https://aclanthology.org/2022.emnlp-main.278>.
- 539 [69] Chenggang Zhao, Shangyan Zhou, Liyue Zhang, Chengqi Deng, Zhean Xu, Yuxuan Liu, Kuai  
540 Yu, Jiashi Li, and Liang Zhao. DeepEP: an efficient expert-parallel communication library.  
541 <https://github.com/deepseek-ai/DeepEP>, 2025.
- 542 [70] Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Y. Zhao, Andrew M.  
543 Dai, Zhifeng Chen, Quoc V. Le, and James Laudon. Mixture-of-experts with expert choice  
544 routing. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh,  
545 editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural  
546 Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28  
547 - December 9, 2022*, 2022. URL [http://papers.nips.cc/paper\\_files/paper/2022/  
548 hash/2f00ecd787b432c1d36f3de9800728eb-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/2f00ecd787b432c1d36f3de9800728eb-Abstract-Conference.html).
- 549 [71] Tong Zhu, Xiaoye Qu, Daize Dong, Jiacheng Ruan, Jingqi Tong, Conghui He, and Yu Cheng.  
550 Llama-moe: Building mixture-of-experts from llama with continual pre-training. *CoRR*,  
551 abs/2406.16554, 2024. URL <https://doi.org/10.48550/arXiv.2406.16554>.
- 552 [72] Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer,  
553 and William Fedus. St-moe: Designing stable and transferable sparse expert models. *CoRR*,  
554 2022. URL <https://arxiv.org/abs/2202.08906>.

555	<b>Contents</b>	
556	<b>1 Introduction</b>	<b>1</b>
557	<b>2 Background and Related Work</b>	<b>2</b>
558	<b>3 Empirical Study</b>	<b>2</b>
559	<b>4 Results</b>	<b>2</b>
560	<b>5 Conclusion</b>	<b>4</b>
561	<b>A Extended Background</b>	<b>13</b>
562	A.1 Continual Pre-training of LLMs . . . . .	14
563	A.2 Mixture of Experts Transformer Language Models . . . . .	14
564	A.3 Training details . . . . .	15
565	A.4 Maximum Routing Imbalance: A proxy for latency in MoEs . . . . .	15
566	<b>B Extended Related Work</b>	<b>15</b>
567	B.1 Mixture of Experts Language Models . . . . .	15
568	B.2 Cotinual Pre-training of Dense Foundation Models . . . . .	16
569	B.3 Continual Pre-training of MoE LLMs. . . . .	16
570	<b>C Extended Experimental Results</b>	<b>17</b>
571	C.1 Ablating Replay (%) and the checkpoint used for CPT . . . . .	17
572	C.2 Language Modeling Performance . . . . .	18
573	C.3 Analyzing changes in routing behaviour due to CPT . . . . .	20
574	C.4 Language Model Evaluation Benchmarks . . . . .	22
575	C.5 Training and Validation Loss . . . . .	24
576	C.6 Qualitative Analysis . . . . .	27
577	C.6.1 Continual routing saturation analysis . . . . .	27
578	C.6.2 Continual Vocabulary Specialization Analysis . . . . .	30
579	C.6.3 Continual expert co-activation analysis . . . . .	34
580	C.6.4 Continual routing imbalance analysis . . . . .	38
581	C.7 Maximum routing imbalance of MoEs During Continual Pretraining. . . . .	41
582	C.8 Limitations . . . . .	41
583	<b>D Dataset Sizes and Sampling Proportions</b>	<b>46</b>
584	<b>E Model Hyperparameters</b>	<b>47</b>
585	<b>A Extended Background</b>	
586	The following section is complementary to section 2 of the main manuscript, providing additional	
587	background for our paper.	



## A.1 Continual Pre-training of LLMs

Continual pre-training extends pre-training to multiple new distributions. Concretely, continual pre-training occurs when a model is trained on a sequence of datasets  $\mathcal{D}_0, \mathcal{D}_1, \dots, \mathcal{D}_N$  with different distributions,  $N \geq 2$ , and each dataset is sufficiently large (e.g.,  $> 100\text{B}$  tokens in the case of language) [26]. Note that the large data scale, here, distinguishes this setting from supervised fine-tuning or instruction tuning where the amount of data is much smaller. Typical application settings of continual pre-training are adapting existing pre-trained models on newly available data or enhancing their capabilities in a specific domain. We will now discuss well-established techniques for continually pre-training dense transformers.

**LR Re-warming and Re-decaying.** Many open-source LLMs follow a linear warmup and cosine annealing schedule during pre-training, which reaches a large maximum learning rate,  $\eta_{\max}$ , early on in training and subsequently decays the learning rate to a small minimum value,  $\eta_{\min}$  [25, 42, 51]. Naively continuing training at  $\eta_{\min}$  or  $\eta_{\max}$  either leads to too little adaptation or too much forgetting. Instead, [26] show that Re-warming and Re-decaying the learning rate during CPT is critical for strong continual learning performance.

**Infinite LR schedules.** While Re-warming and Re-decaying the learning rate following a cosine decay schedule was found to be a good solution when starting from a fully decayed checkpoint, [26] remark that this strategy incurs forgetting, due to the large LR increase, even when continually pre-training on the same distribution. To circumvent this, the authors propose infinite learning rate schedules that allow for a smooth transition in learning rate between continual learning phases and are not bound to a fixed number of training steps.

**Replay.** Replaying previous data has been a longstanding tool for mitigating catastrophic forgetting [64]. In our experiments, we replay previously seen data for this purpose, designating any model using the technique with the suffix “ $X\%$  Replay”. Here,  $X$  represents the percentage of samples in a given batch that were replayed from the previous distribution. To match compute across different replay budgets, we do not increase the token budget when increasing the amount of replay. Instead, we decrease the amount of new data seen during CPT.

## A.2 Mixture of Experts Transformer Language Models

Sparsely-activated MoE transformers differ from their dense counterparts by dynamically routing tokens in a sequence,  $X \in \mathbb{R}^{S \times H}$ , to different experts  $\{\text{FFN}_{i,j}(\cdot)\}_{i=0}^N$  as opposed to a single FFN. Here  $S$  is the sequence length,  $H$  is the transformer’s hidden dimension,  $j$  indexes over the transformer’s blocks, and  $N$  is the number of experts per block. This is often referred to as an MoE layer [58]. Typically, these layers are used in place of Feed Forward Networks (FFN) in each transformer block [18, 12], however, recent works [59, 68] have also replaced the query and output matrices of multi-head self-attention layers with MoE layers. In what follows, we exclusively study MoE transformers that replace FFNs at each block with MoE layers, similarly to state-of-the-art recent work [12, 62, 46, 16, 15]. Moreover, we also study the recent trend of using more granular experts and shared experts [12, 62, 46, 24, 40, 43, 53, 16, 15].

Algorithms for dynamically selecting among experts, known as routing algorithms [55, 58, 72, 8, 37], are central to MoEs. A key consideration for token-choice (we do not consider expert-choice as it is incompatible with autoregressive generation) routing algorithms is achieving a balanced load across experts in a given layer. Without enforcing a balanced load, the router may collapse to only choosing a single or a few experts, leading to poor parameter utilization and higher latency proportional to the load of the most burdened expert [70].

In this work, we focus on studying two prominent Top- $k$  routing algorithms from recent state-of-the-art works which we refer to as Penalty-Balanced Top- $k$  (PBT $k$ ) routing [58, 12, 72, 18] and Sinkhorn-Balanced Top- $k$  (SBT $k$ ) routing [8, 4]. Both algorithms define the router  $R(x) = Wx : \mathbb{R}^H \rightarrow \mathbb{R}^e$  to be a simple linear projection to the space of experts. Expert probabilities are computed by applying a softmax to the router’s output:  $p(x) = \text{softmax}(SB(R(x)))$ . Where  $SB(\cdot)$  is the Sinkhorn load balancing function in the case of SBT $k$  routing or the identity otherwise.

Ranked by  $p(x)$ , the top- $k$  experts are selected for each token, where  $k$  is a hyperparameter selected before training. For a single token, the output of MoE layer  $j$ ,  $f_{\text{MoE}_j}$ , is computed as follows:

$$f_{\text{MoE}_j}(x) = \text{SFFN}_j(x) + \frac{\sum_{i \in I_k(x)} p_i(x) \cdot \text{FFN}_{i,j}(x)}{\sum_{i \in I_k(x)} p_i(x)}.$$

Where  $I_k(x) = \{i \mid p_i(x) \in \text{Top } k \text{ elements of } p(x)\}$  and SFFN is a shared expert [53, 12] if one is used or the identity otherwise. While they both route tokens to the top- $k$  experts, the PBT $k$  and SBT $k$  routing algorithms differ in how they balance the load across experts.

**Penalty-Balanced Top- $k$  Routing.** PBT $k$  methods in the literature [58, 18, 72, 12] add penalty terms to the overall loss to encourage a balanced load across experts. The *auxiliary loss* has become the most popular such penalty and is used in conjunction with the  $z$ -loss in several recent state-of-the-art MoEs [12, 62]. Briefly, the auxiliary loss is minimized when the router assigns an equal proportion of tokens in a given batch to each expert in a given block, while the  $z$ -loss penalizes large-magnitude router logits. The latter has been shown to promote numerical stability in larger models [72]. Given their combination in recent SOTA MoE LLMs [72, 12], we exclusively study MoEs that combine both *auxiliary loss* and  $z$ -loss, referring to them as PBT $k$  MoEs.

**Sinkhorn-Balanced Top- $k$  Routing.** SBT $k$  routing casts the assignment of tokens to experts as a linear assignment problem which corresponds to a well-studied problem in optimal transport, namely "the regularized Kantorovich problem of optimal transport" [8]. The Sinkhorn-knopp algorithm [33] provides an approximate solution to this problem which can be efficiently computed on GPUs. In practice, this corresponds to adjusting routing probabilities (e.g. according to  $SB(\cdot)$ , see [8] section B.2.1 for details) such that a relatively balanced load is obtained without deviating *too much* from greedy Top- $k$  routing.

### A.3 Training details

All models in our study (except re-training baselines) were pre-trained for 192, 720 gradient descent steps using a batch size of 1024, a sequence length of 2048, the AdamW optimizer, and the *Cosine Inf* schedule [26]. For continual pre-training, each model in the main study follows a *Cosine Inf* schedule resuming from the non-decayed checkpoint, while some models in the ablation section were continually pre-trained from decayed checkpoints following a cosine decay schedule (e.g., replicating the setting from [26]). We continually pre-train the models for 95, 370 gradient descent steps using the same batch size and sequence length as during pre-training. Each model was trained across 64 A100 GPUs using data parallelism and zero-1 [52]. To accelerate the droppless MoE forward pass we use the Megablocks kernel [20].

### A.4 Maximum Routing Imbalance: A proxy for latency in MoEs

While performance is one important axis of robustness to distribution shifts, maintaining a balanced load across experts is just as important for MoE foundation models. Without a balanced load, MoE transformers inferred using expert parallelism without token dropping (e.g., as is done for SOTA models [16, 69]) could be bottlenecked by the speed of a single accelerator that receives all the tokens, leading to underutilization of the hardware, lower throughput, and higher costs. To quantitatively assess the effect of distribution shift on load balance, we propose the maximum routing imbalance (MRI): the largest proportion of tokens routed to a single expert in a given MoE layer. Concretely, the MRI at a training iteration  $t$  and MoE layer  $j$  is defined as

$$\text{MRI}(t, j) := \max_{i \in [1, \dots, E]} \left[ \frac{\sum_{x \in B} \mathbb{1}\{i \in I_k(x)\}}{|B|} \right]. \quad (1)$$

Where  $B$  is a set containing all tokens in a given batch,  $\mathbb{1}$  is the indicator function,  $E$  is the number of routed experts, and  $k$  is the number of active experts. *Since latency increases with computation, and, in an MoE layer, the computation required by a given device increases with the load of experts on that device, then MRI calculated with respect to routing decisions on a distribution is a proxy for the worst case latency of an MoE layer on the distribution.* We will use the MRI throughout the following sections to measure the effect of algorithmic changes to continual pre-training on routing imbalance.

## B Extended Related Work

### B.1 Mixture of Experts Language Models

Mixture of experts language models have a long history with fundamental ideas dating back several decades [10, 27]. More recently, in the context of large-scale language modeling, the mixture-of-experts layer [58] was introduced to substantially increase the capacity of an LSTM language model

with little detriment to efficiency. The authors also introduced a load-balancing penalty to encourage even utilization of experts. Subsequently, Fedus et al. [18] refined the penalty, renamed *auxiliary loss*, which has become a central component of modern MoEs. Follow-up works have focused on massively scaling up MoE LLMs, improving the routing algorithms of these models, improving the quality of the router’s gradient estimate, and making architectural improvements to these models. Lepikhin et al. [36] introduce the MoE layer into the transformer architecture, using two activated experts (thought to be necessary for nontrivial router gradients) and scaling to an unprecedented 600B parameter scale. Subsequently, Fedus et al. [18] introduced the Switch Transformer, showing that it is possible to scale MoEs beyond 1T parameters despite training them with only a single active expert.

Other works have focused on developing novel routing algorithms. Lewis et al. [37] cast routing as a linear assignment problem and leverages Hungarian matching in their routing algorithm. Clark et al. [8] use Sinkhorn’s algorithm to approximately solve the assignment problem on GPUs, resulting in a faster algorithm. Anthony et al. [4] introduce a favorable initial condition to improve the convergence of the iterative Sinkhorn solver, further reducing the cost of sinkhorn routing. Roller et al. [55] introduces a deterministic routing algorithm based on hash layers. Zoph et al. [72] introduces a loss penalty to promote stability in large-scale MoE routing. Wang et al. [63] introduces the first learned routing mechanism that uses neither an entropy regularizer nor an assignment-based approach to balance expert utilization in token-choice routing. [70] introduce Expert Choice Routing, a routing paradigm where each expert receives a balanced load and the routing algorithm decides which tokens to send to each of the experts; while it obtains strong performance and automatically achieves a balanced load, ECR is incompatible with autoregressive generation so we don’t consider it in this work. Other works propose methods to better approximate the full router gradient [47, 39, 38].

Finally, a recent trend of using MoE experts with finer-grained intermediate sizes has shown notable performance gains when compared to using the full intermediate FFN size, as was originally done [58, 18, 36]. Liu et al. [40] first observed that utilizing smaller expert layers improves perplexity. Subsequently, researchers have explored scaling laws for fine-grained MoEs at small scale [43], pre-trained and released SOTA MoEs that leverage the fine-grained expert architecture [12, 62, 46], and pushed the idea of thinner experts to its limit, exploring MoEs with millions of experts [24]. While we have reviewed the most relevant works to ours, there are many more works that we have not had the chance to mention here. We refer the avid reader to a recent and comprehensive survey of the area [6].

## B.2 Cotinual Pre-training of Dense Foundation Models

Continual pre-training of foundation models has the same objectives as continual learning [19], except that it is applied to large-scale pre-training tasks, which are mainly self-supervised. Several existing works study continual learning in settings relevant to continual pre-training. They find that self-supervised pre-training benefits from reduced forgetting [11, 13], that pre-trained models forget less than their randomly initialized counterparts [44], that forgetting improves as model scale is increased [54], and that wider models tend to forget less than deeper models [45]. In the context of LLM fine-tuning, [57] shows that little replay is needed to prevent forgetting when fine-tuning on small amounts of instruction-tuning data. In the context of large-scale (with respect to data) continual pre-training for LLMs, Gupta et al. [23] highlights the importance of rewarming the learning rate when models are pre-trained from a checkpoint that has been decayed to a small learning rate. Following up on their work, Ibrahim et al. [26] establish the effectiveness of learning rate re-warming, LR re-decaying, and replay for large-scale continual pre-training of LLMs. Concurrently, Garg et al. [21] establishes the performance of the same techniques for CLIP models. Shortly thereafter, Parmar et al. [48] scale continual pre-training for dense decoder-only transformers further, showing that a 15B parameter model pre-trained for 8T tokens can be effectively pre-trained on 1T tokens of incoming data.

## B.3 Continual Pre-training of MoE LLMs.

To the best of our knowledge, only a single work exists exploring the large-scale continual pre-training of MoEs LLMs, while the majority of the literature focuses on upcycling or growing MoEs for continual pre-training.

In a concurrent pre-print DeepSeek-CoderV2 [14], shows that they can continue from a checkpoint the training of a MoE LLM. However, this is only shown for one instance and the analysis of the MoE routing behavior is not discussed. Furthermore, there is no comparison to a FLOP-matched dense model, making it challenging to assess whether the sample efficiency of MoE LLMs is maintained during continual pre-training.

Continual pre-training methods for MoEs that are less related to our work generally focus on fine-tuning MoE LLMs on small amounts of data [65] or growing MoEs [35, 71, 61, 22]. Wang et al. [65] study MoE specific techniques for parameter-efficient fine-tuning (PEFT). Zhu et al. [71] proposes a technique to create an MoE by splitting the FFNs of an existing dense transformer and subsequently continually pre-training it. Sukhbaatar et al. [61] proposes to continually pre-train a dense LLM on multiple different datasets, gather the FFN layers from different continually pre-trained models to form MoE layers, merge the parameter tensors other than FFN layers, and subsequently continually pre-train the merged model to learn routing in the MoE part. Gritsch et al. [22] propose a similar method to train new expert layers that uses domain embeddings from a pre-trained embedding model as the identifier for a domain’s experts, allowing the domain embeddings to provide an inductive bias that can help with adding new experts. While these methods allow for improving the capabilities of MoEs with new data, they focus on first upcycling dense models, whereas we focus on updating MoEs pre-trained from scratch.

## C Extended Experimental Results

In the following section, we provide extended experimental results from the paper in a non-summarized format to enhance the reproducibility of our manuscript and allow the reader to dive into whichever details may most interest them.

### C.1 Ablating Replay (%) and the checkpoint used for CPT

In the following section, we ablate the replay percentage used during continual pre-training and consider continually pre-training from two distinct checkpoints: a checkpoint that (a) was decayed to  $\eta_{min}$  during pre-training (the case for most open-source MoEs) or (b) followed an infinite learning rate schedule and was not decayed (the ideal case achievable when one has control over the pre-training phase). Models in group (a) are continually pre-trained following a linear warmup and cosine decay schedule that rewarms the learning rate to  $\eta_{max}$  before re-decaying it (e.g., as in [26]), while the models in group (b) are continually pre-trained starting from  $\eta_{const}$  following an infinite LR schedule (exact values are provided in Sec. E).

**Validation Loss.** Figure 2 reports the validation loss for these models during the first 50B tokens of continual pre-training. While we only show the first 50B tokens due to resource constraints, the schedules were set to decay at 200B tokens, mimicking the start of a longer continual pre-training. Subfigures (a) and (c) show forgetting and adaptation plots using 0% replay, while subfigures (b) and (d) show analogous plots using 40% replay. We observe that as the percentage of replay is increased, the forgetting as measured by FineWeb validation loss is mitigated, while the adaptation to the downstream dataset is harmed. Turning our attention to the checkpoints used, we observe that, for all replay values and all models, using non-decayed checkpoints improves forgetting on the FineWeb without compromising adaptation. These results show that, similar to dense transformers, MoEs can tradeoff forgetting for adaptation with replay and benefit from infinite LR schedules

**Routing Imbalance.** Figure 4 (a,b) reports median MRI computed across all transformer blocks of SBT $k$  and PBT $k$  Granular MoEs during CPT with 0% replay, subfigure (c) reports results across different replay percentages and results for switch MoEs are reported in Figures 19 and 18 of the appendix. These figures precisely study the distribution shift by reporting the MRI immediately before and after the transition from English and German data. We observe that SBT $k$  MoEs are consistently robust to the distribution shift, showing only a small increase in MRI across different replay percentages and for decayed and non-decayed checkpoints alike. This is likely attributable to the explicit balancing step in Sinkhorn routing. In contrast, the non-decayed and decayed PBT $k$  MoE checkpoints go through a period of high routing imbalance immediately following the distribution shift. However, this period is short-lived: the PBT $k$  checkpoints recover to well-balanced MRI levels, better than those of SBT $k$ , within 500 training steps. Subfigure (c) shows that this MRI spike can be mitigated with replay, though the benefit is negligible, as even the no-replay model recovers quickly.

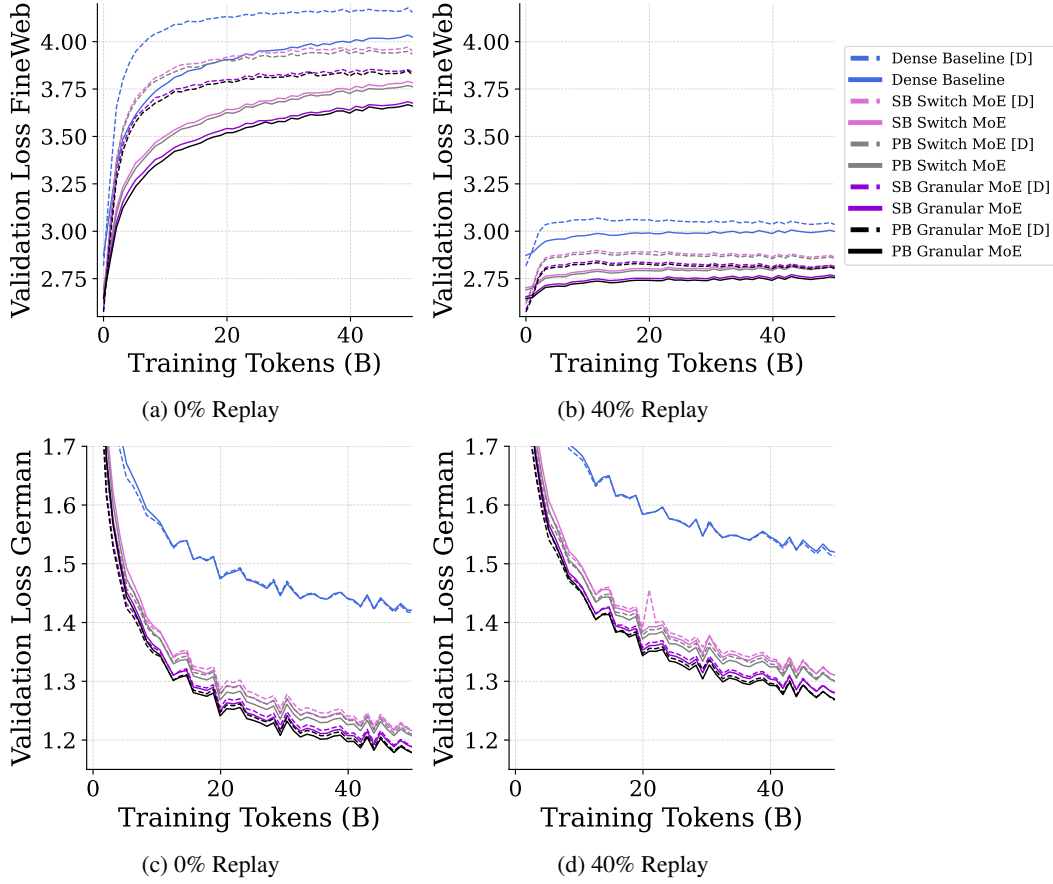


Figure 2: **Ablating replay and decay strategy during continual pre-training on German data.** We CPT MoEs and a dense baseline from fully-decayed checkpoints (dotted curves, [D]) or a non-decayed checkpoint (full curves). The figures report the performance on task 1 (FineWeb) and task 2 (German) while CPT on task 2. We observe that adaptation to task 2 is similar within an architecture for both checkpoints, that CPT from a non-decayed checkpoint improves forgetting, and that replay mitigates forgetting.

795 These results suggest that although SB is more robust to distribution shifts than PB, this robustness  
796 limits the MRI attainable and could be the cause of the poorer performance seen above for validation  
797 loss. Moreover, the chaotic phase undergone by PBT $k$  checkpoints does not last long enough to  
798 forego the strong performance of these models.

## 799 C.2 Language Modeling Performance

800 Having established the benefits of replay and infinite learning rate schedules for continually pre-  
801 training MoEs, we now quantitatively verify the efficacy of these techniques by continually pre-  
802 training our MoEs on 200B tokens of Code and German Common Crawl data and evaluating their  
803 performance relative to two baselines. Specifically, we will compare the performance of the four  
804 continually pre-trained MoEs in our study to a FLOP-matched dense baseline and a fully re-trained  
805 PBT $k$  Granular MoE Baseline (the best-performing architecture in our study). Performance will be  
806 measured across 4 axes: validation loss on the pre-training and CPT datasets, English evaluation  
807 benchmarks (task 1), German and Code evaluation benchmarks (task 2), and MRI of final checkpoints.  
808 Note that the main conclusions of this section are succinctly summarized in Figure 3.

809 **Validation Loss.** Table 1 reports validation loss (e.g., log perplexity on a fixed held-out validation set)  
810 results for the main models in our study, while extended results are reported in Table 5 of the appendix.  
811 We observe that all MoEs outperform the FLOP-matched dense baseline during pre-training and  
812 CPT. Within the MoEs, we observe that PBT $k$  MoEs consistently outperform SBT $k$  MoEs and that



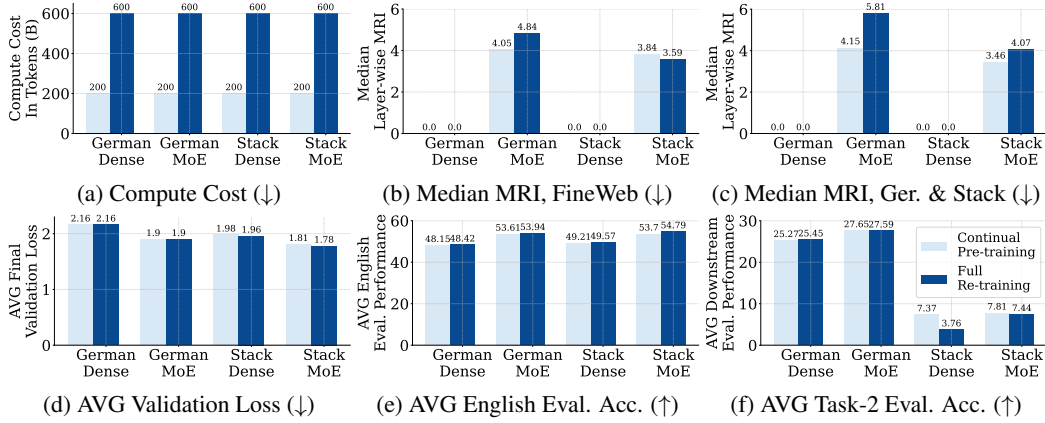


Figure 3: **Continually pre-trained MoEs match the performance of full re-training across two dataset transitions: 400B English → 200B German and 400B English → 200B Stack.** We compare the performance of a fully re-trained Penalty-Balanced Top- $k$  MoE and dense baseline, to their CPT counterparts. Despite incurring only a third of the substantial full-retraining cost, the CPT MoEs match the performance of the fully re-trained models, even achieving improvements in median Maximum Routing Imbalance (MRI) in some cases. This shows that MoEs have CPT abilities on par with dense transformers. Note that subfigures (b), (c), and (f) evaluate German and Stack models on different datasets which correspond to their training domain.

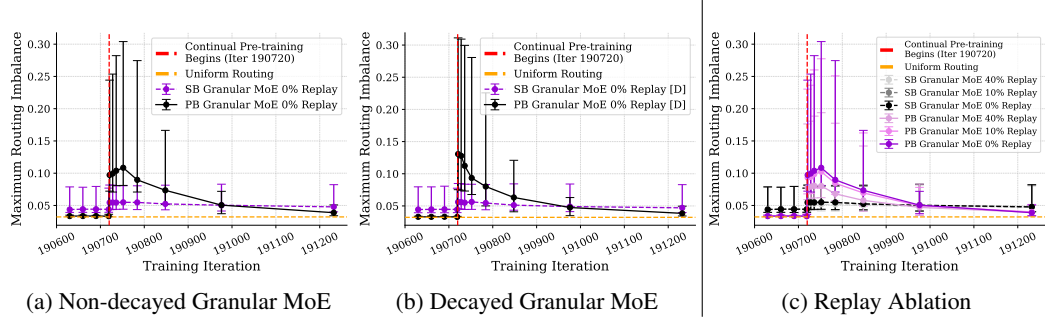


Figure 4: **FineWeb → German CPT checkpoint and replay ablation.** We report the median Maximum Routing Imbalance (MRI) across MoE layers with min/max error bars. Sinkhorn-Balanced (SBT $k$ ) MoEs show a slight MRI increase during distribution shift, while PBT $k$  MoEs experience a brief spike before recovering to balanced MRI levels below SBT $k$ , which approach the uniform routing baseline. The uniform routing baseline (orange line) corresponds to the case where each expert across all layers receives the same number of tokens; thus, it represents perfect balance.

Granular MoEs consistently outperform switch MoEs across pre-training and CPT. These findings are consistent with the literature on Granular MoEs [43, 12], but we believe this is the first time that SBT $k$  routing has been shown to underperform PBT $k$  routing. Since the PBT $k$  Granular MoE achieves the best performance, we use it as our full re-training baseline. Compared to full re-training, our Granular CPT MoEs consistently have higher FineWeb validation loss but achieve lower downstream validation loss with a similar average validation loss. Similar results are found when comparing the CPT dense baseline to its full re-training counterpart. These results demonstrate that the continual learning abilities of MoEs w.r.t. validation loss are on par with dense models for adaptation and are slightly superior in terms of forgetting, likely due to their larger total parameter count.

**English Evaluation results.** Table 1 presents average accuracy, while Table 3 details per-benchmark results. We select benchmarks where models of our scale (570M active parameters) achieve non-trivial accuracy to maximize signal. Each model is evaluated zero-shot on benchmarks covering Commonsense Reasoning, Reading Comprehension, Scientific Question Answering, and Math: HellaSwag [67], Winogrande [56], PIQA [5], ARC-Easy, ARC-Challenge [9], SWAG [66], LAMBADA (OpenAI) [60], SciQ [31], PubMedQA [30], and MathQA [3] (see Sec. C.4 for details). We find that models trained solely on FineWeb outperform all others, including full re-training baselines. Granular MoEs surpass switch MoEs. CPT models trained on Stack perform similarly to those

trained on German. Compared to full re-training, CPT models achieve nearly identical results (within  $\sim 1\%$ ). The dense baseline also matches its full re-training counterpart, indicating that MoEs have similar continual learning abilities on pre-training evaluations while benefiting from improved sample efficiency.

**German Evaluation results.** Table 1 shows average German evaluation performance, while table 4 of the appendix provides a per-benchmark breakdown. We use GPT-3–translated German versions of HellaSwag, ARC-Challenge, and TruthfulQA, evaluating each zero-shot [50]. German-trained models outperform English-only ones, and German-trained MoEs surpass the FLOP-matched dense baseline. Among MoEs, modules using the same training tokens perform similarly. CPT MoEs and the full re-training baseline differ by  $< 1\%$  accuracy, with no clear winner. The dense baseline also performs comparably to full re-training, demonstrating that the continual learning abilities of MoEs w.r.t. German evaluations are on par with dense models while benefiting from improved sample efficiency.

**Code Evaluation results.** Table 1 presents average Code evaluation performance, while Table 2 of the appendix provides a pass@ $k$  breakdown ( $k \in \{1, 10, 50, 100, 150, 200\}$ ). Our models are evaluated on Python code-generation tasks from HumanEval [7], as Python is well-represented in our Stack CPT dataset (Table 9). English-trained models can not solve any problem, whereas stack-trained models achieve non-trivial accuracy. Unlike for other performance metrics, CPT switch MoEs slightly outperform their granular counterparts. Compared to full re-training, all CPT MoEs perform marginally better, while the CPT dense model exceeds its baseline by over 3%. We attribute this unexpected improvement to evaluation noise and training variance, given the models’ similar validation loss. These results suggest MoEs match dense models in continual learning for code evaluation when accounting for MoEs’ improved sample efficiency.

**Routing imbalance during and after continual pre-training.** Figure 1 shows the layer-wise Maximum Routing Imbalance (MRI) for Granular MoEs across FineWeb (a), German (b), and stack (c), while Figure 16 reports MRI for all MoEs. We include a 0% replay baseline for each MoE CPT on German to highlight replay’s impact on MRI.

In subfigure (a), Penalty-Balanced MoEs consistently have lower MRI than Sinkhorn-Balanced MoEs across all architectures, and granular MoEs exhibit lower and more stable MRI within architectures. On FineWeb, continual pre-training causes only a slight MRI increase relative to the pre-trained checkpoint, even for the 0% replay model, except in its first layer. Interestingly, all German-trained MoEs show higher MRI on FineWeb than their Stack-trained counterparts, with the full re-training baseline, surprisingly, having the highest. This suggests there may be more routing interference between English and German datasets and that continual pre-training may help reduce MRI across distributions, possibly due to the use of *CosineInf* vs. *Cosine Annealing* schedules for CPT and re-training, respectively.

Granular MoEs also reduce routing imbalance on German (b) and Stack (c) datasets. MoEs become most unbalanced with out-of-distribution data (e.g., non-German models in (b) and non-code models in (c)). Similar trends hold for Switch MoEs, with an additional finding: high MRI is common in early layers of Switch MoEs, independent of training/testing distributions, unlike Granular MoEs. These results show that PBT $k$  and SBT $k$  MoEs are robust to distribution shifts w.r.t. MRI and can even outperform re-training baselines, suggesting that continually pre-training MoEs should have no negative impact on inference latency.

*In summary, we find that, across three measures of performance, MoEs continually pre-trained with replay and infinite LR schedules can match the performance of a full re-training baseline and, thus, they have similar CPT abilities to a FLOP-matched dense baseline without any inhibition from their routers. Moreover, we show that continually pre-training MoEs has no negative impact on MRI compared to re-training.*

### C.3 Analyzing changes in routing behaviour due to CPT

In the following section, we analyze changes in routing behaviour resulting from continual pre-training. Specifically, we record routing decisions of the MoE checkpoints before and after continual pre-training on 20,000,000 tokens of held-out test data from FineWeb, Stack, and German. To understand how routing decisions change during CPT, we adapt three routing behavior metrics from [46] to the continual pre-training setting: Router Saturation, Vocabulary specialization, and

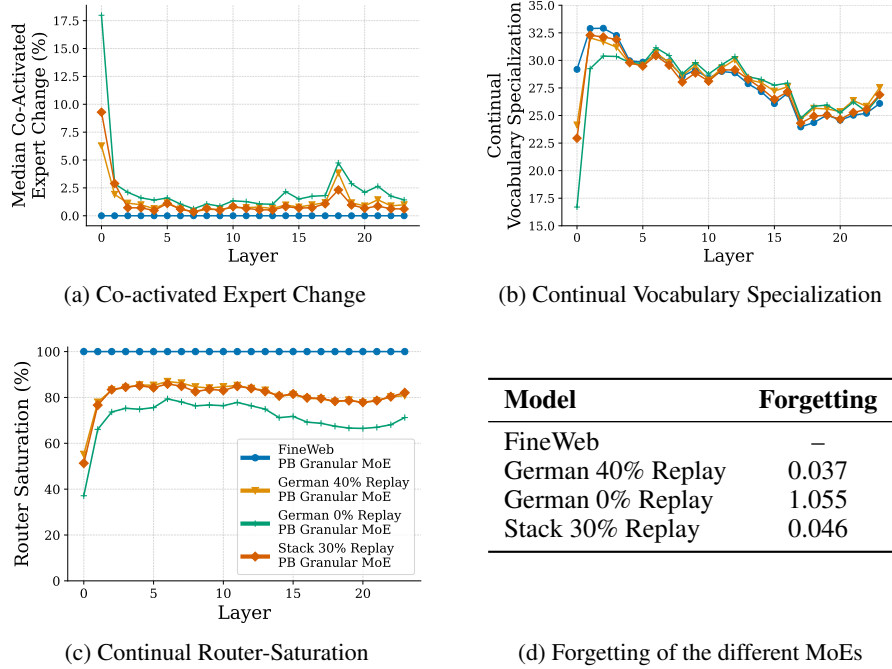


Figure 5: **Layer-wise analysis of routing changes during CPT.** Our goal is to understand how routing decisions change from the pre-trained checkpoints to final checkpoints after continual pre-training. To this end, we analyse changes in routing behaviour from 3 perspectives: which experts tend to be activated together (a), the tendency for certain vocabulary tokens to be routed to certain experts (b), and how close routing decisions of the pre-trained checkpoint are from CPT checkpoints (c). To provide context to these metrics, we remind the reader of the forgetting (e.g., from table 1) for each model shown in the plots. We observe that the no-replay baseline changes the most in early layers and forgets the most, suggesting that more drastic changes in initial layers may be linked to forgetting.

884 Expert co-activation. We will provide brief descriptions of each in what follows, with formal  
885 definitions available in the appendix (Sections C.6.1, C.6.2, and C.6.3).

886 **Continual Router-Saturation** Router Saturation (RS) is the percentage of routing decisions at iteration  $t$   
887 that match those of the final checkpoint [46]. We extend this metric to multiple training phases  
888 for continual pre-training. Figure 5 (c) shows RS between the pre-training and CPT checkpoints for  
889 Stack and German Granular PBTk MoEs. RS is lowest in early layers, peaks at layers 2 – 13, and  
890 slightly drops after layer 13. The 0% replay German checkpoint has RS consistently 10 – 15% lower  
891 than the 40% replay checkpoint across all layers. Note that despite adapting well to German, only  
892 the no-replay checkpoint suffers significant forgetting on FineWeb. These results suggest that CPT  
893 adaptation is most pronounced in layers 0 – 2 and 13 – 23 and that forgetting has the same pattern  
894 but is correlated with lower overall router saturation.

895 **Continual Vocabulary Specialization** Vocabulary Specialization (VS) quantifies how often a token  
896 from a dataset is routed to a specific MoE expert relative to its total occurrences [46]. By assigning  
897 each token in the model’s vocabulary to the expert that processes it the most frequently, we can  
898 create a *one-to-many mapping* between experts and vocabulary entries for MoE layer. Then, we can  
899 calculate the average VS of each expert by averaging over its assigned tokens and averaging across  
900 experts to measure specialization within a layer. To compare specialization across model checkpoints,  
901 we can re-use the *one-to-many mapping* of a previous checkpoint and measure how the specialization  
902 w.r.t. this mapping has changed during CPT. Figure 5 shows VS w.r.t. pre-training checkpoints using  
903 FineWeb data. VS is notably lower in layers 0-4 after CPT while there is almost no discernable  
904 change in VS for layers 5 – 23. The zero-replay checkpoint exhibits the lowest VS in layers 0-4,  
905 correlating with its weaker FineWeb performance and suggesting that excessive VS shifts in early  
906 layers may contribute to forgetting.

**Co-activated Expert Change** Expert co-activation between two experts  $E_i$  and  $E_j$  is defined as the ratio of times they are activated together to the total activations of  $E_i$  over a dataset [46]. This metric applies only to MoEs with  $k \geq 2$  active experts. A co-activation matrix can be constructed for each ordered expert pair in a layer. To compare expert co-activation across model checkpoints, we compute co-activation matrices for all layers of two checkpoints ( $C^{(1)}, C^{(2)}$ ) and measure absolute changes by computing statistics of the entries in their element-wise absolute differences matrix ( $|C^{(1)} - C^{(2)}|$ ). Figure 5 (a) shows the median co-activation change between the pre-training and CPT checkpoints. Early layers (0-1) exhibit the largest changes, with a consistent spike at layer 18 for all CPT models and slightly higher median changes in layers 13 – 23. Among CPT checkpoints, the no-replay variant shows the most significant co-activation shifts. Despite all checkpoints adapting well to new distributions, only the no-replay checkpoint experiences substantial forgetting on FineWeb. These findings suggest that adaptation during CPT correlates with co-activation changes in early (0 – 2) and later (13 – 23) MoE layers and that more pronounced changes correlate with higher forgetting.

#### C.4 Language Model Evaluation Benchmarks

We evaluate the language models in our study on English, Code, and German evaluation tasks. Please note that our goal is not to achieve SOTA performance on these benchmarks; none of our models have been aligned or fine-tuned to improve performance. Instead, we seek to evaluate their performance within the context of our controlled scientific study. Given the scale of our language models (at most 570M active parameters), we carefully select evaluation tasks that show non-trivial evaluation results; that is, we choose tasks for which the models in our suite achieve above random chance accuracy.

Selected English evaluation tasks:

- **Commonsense Reasoning (0-shot):** HellaSwag [67], Winogrande [56], PIQA [5], ARC-Easy, ARC-Challenge [9], SWAG [66]
- **Reading Comprehension (0-shot):** LAMBADA (OpenAI) [60]
- **Scientific Question Answering (0-shot):** SciQ [31], PubMedQA [30]
- **Math (0-shot):** MathQA [3]

Selected German evaluation tasks translated from the corresponding English language tasks using the GPT 3.5 API [50].

- **Commonsense Reasoning (0-shot):** HellaSwag-DE [67], ARC-Challenge-DE [9]
- **Reading Comprehension (0-shot):** TruthfulQA-DE [32]

Code evaluation tasks

- **Python:** Human Eval (pass@1-200)

Tables 3, 2, and 4 report the performance of models in our study on English, Code, and German evaluation benchmarks.

Table 2: **Human Eval after pre-training on FineWeb and continual pre-training on Stack.** We report the percentage of problems for which at least one generated solution passes all tests. We observe that all English-only models generate only incorrect solutions, while the models continually pre-trained on code and the full re-training baselines achieve non-trivial accuracy. Interestingly, the SB Switch MoE performs best of all across all pass thresholds. However, given the generally poor performance of the models overall, we attribute differences within a dataset type to random chance.

Training Tokens	Model	pass@1	pass@10	pass@50	pass@100	pass@150	pass@200	Mean
400B FineWeb	Dense Baseline	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
	SB Switch MoE	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
	PB Switch MoE	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
	SB Granular MoE	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
	PB Granular MoE	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
400B FineWeb → 200B Stack	Dense Baseline	0.21%	1.94%	6.87%	9.96%	11.85%	13.41%	7.37%
	SB Switch MoE	<b>0.24%</b>	<b>2.19%</b>	<b>8.06%</b>	<b>12.15%</b>	<b>14.85%</b>	<b>17.07%</b>	<b>9.09%</b>
	PB Switch MoE	0.21%	1.93%	7.28%	11.10%	13.56%	15.24%	8.22%
	SB Granular MoE	0.18%	1.69%	6.50%	10.01%	12.30%	14.02%	7.45%
	PB Granular MoE	0.16%	1.51%	6.30%	10.40%	13.24%	15.24%	7.81%
400B FineWeb ∪ 200B Stack	Dense Baseline	0.14%	1.20%	3.57%	4.99%	5.98%	6.71%	3.76%
	PB Granular MoE	0.20%	1.82%	6.84%	10.21%	12.13%	13.41%	7.44%

Table 3: **Language models evaluation benchmarks after pre-training (English Web Data), continual pre-training (Code & German Web Data), and full re-training.** We report accuracy for all selected benchmarks. We observe that all MoEs and the dense baselines maintain similar relative performance before and after the distribution shift, showing that MoE LLMs’ continual pre-training dynamics are similar to dense models with respect to forgetting on evaluation tasks. When comparing the mean evaluation performance of PB granular MoE to the full re-training baseline, we observe that the final performance is nearly reached or matched with a substantially smaller computational cost.

Training Tokens	Model	ARC-C	ARC-E	HellaSwag	LAMBADA OAI	MathQA	PIQA	PubMedQA	SciQ	SWAG	WinoGrande	Mean
400B FineWeb (Annealed)	Dense Baseline	23.55%	54.25%	39.99%	47.55%	23.52%	71.98%	51.80%	82.70%	47.59%	55.49%	49.84%
	SB Switch MoE	26.79%	60.48%	45.60%	53.44%	24.52%	74.16%	<b>62.00%</b>	84.80%	50.31%	59.27%	54.14%
	PB Switch MoE	<b>28.75%</b>	61.15%	46.16%	54.22%	<b>25.86%</b>	74.37%	58.20%	87.20%	50.67%	57.93%	54.45%
	SB Granular MoE	26.19%	<b>65.19%</b>	48.10%	56.24%	24.66%	74.92%	61.10%	89.30%	51.35%	60.06%	<b>55.71%</b>
	PB Granular MoE	<b>28.75%</b>	62.92%	<b>48.45%</b>	56.08%	24.62%	75.24%	60.00%	88.90%	<b>51.36%</b>	59.59%	<b>55.59%</b>
400B FineWeb (Non-Annealed)	Dense Baseline	22.35%	52.02%	39.12%	49.04%	23.15%	70.46%	52.50%	81.90%	47.09%	54.14%	49.18%
	SB Switch MoE	24.23%	57.91%	44.26%	51.72%	24.52%	73.12%	<b>60.30%</b>	83.50%	49.32%	56.07%	52.55%
	PB Switch MoE	26.54%	60.86%	44.59%	53.21%	23.99%	73.39%	52.30%	85.80%	49.98%	56.27%	52.69%
	SB Granular MoE	26.54%	62.63%	46.85%	55.87%	24.56%	73.67%	58.60%	87.70%	50.72%	59.04%	54.62%
	PB Granular MoE	27.82%	61.20%	46.52%	55.46%	24.36%	74.97%	58.80%	86.80%	50.87%	58.64%	54.54%
400B FineWeb → 200B German (40% Replay)	Dense Baseline	22.87%	51.43%	36.97%	46.75%	23.75%	70.18%	48.80%	80.70%	45.87%	54.22%	48.15%
	SB Switch MoE	24.66%	56.90%	42.99%	52.94%	24.22%	73.07%	57.40%	83.50%	48.85%	55.41%	51.99%
	PB Switch MoE	25.34%	56.94%	42.59%	53.43%	25.06%	73.23%	49.40%	84.20%	48.76%	53.51%	51.25%
	SB Granular MoE	25.43%	60.02%	44.67%	55.23%	25.13%	73.01%	55.10%	84.60%	49.89%	<b>60.46%</b>	53.35%
	PB Granular MoE	27.05%	60.52%	44.66%	54.43%	24.56%	73.88%	58.40%	85.70%	49.70%	57.22%	53.61%
400B FineWeb ∪ 200B German CC	Dense Baseline	23.29%	51.35%	36.77%	46.17%	24.19%	70.08%	54.00%	80.30%	45.51%	52.57%	48.42%
	PB Granular MoE	27.99%	60.06%	45.06%	55.23%	25.16%	73.50%	56.20%	86.20%	49.88%	60.14%	53.94%
400B FineWeb → 200B Stack (30% Replay)	Dense Baseline	22.01%	52.95%	37.49%	46.98%	22.91%	71.06%	55.40%	83.70%	45.76%	53.83%	49.21%
	SB Switch MoE	22.87%	55.98%	42.51%	52.84%	24.12%	72.80%	55.80%	85.10%	48.78%	56.83%	51.76%
	PB Switch MoE	26.28%	59.01%	42.78%	53.17%	24.32%	73.50%	55.10%	86.20%	49.07%	56.43%	52.59%
	SB Granular MoE	26.54%	60.19%	44.57%	55.44%	24.39%	73.01%	55.60%	85.90%	49.83%	59.59%	53.51%
	PB Granular MoE	25.43%	60.27%	44.88%	54.94%	25.36%	73.78%	56.90%	88.00%	49.62%	57.85%	53.70%
400B FineWeb ∪ 200B Stack	Dense Baseline	22.18%	52.78%	38.68%	48.50%	24.49%	71.00%	51.70%	83.40%	47.01%	55.96%	49.57%
	PB Granular MoE	27.82%	62.67%	46.43%	<b>56.39%</b>	25.66%	<b>75.35%</b>	56.60%	<b>89.40%</b>	50.85%	56.75%	54.79%

Table 4: **German Language models evaluation benchmarks after pre-training (English Web Data), continual pre-training (Code & German Web Data), and full re-training.** We report accuracy for all selected benchmarks. We observe that all MoEs and the dense baseline improve performance on German after continual pre-training. When comparing to the full re-training baselines, we observe that the average performance of our continually pre-trained models are on par.

Training Tokens	Model	Arc-C DE	Hellaswag DE	TruthfulQA DE (MC1)	Mean
400B FineWeb	Dense Baseline	18.52%	26.78%	25.34%	23.54%
	SB Switch MoE	18.60%	26.87%	23.87%	23.11%
	PB Switch MoE	18.94%	26.56%	24.60%	23.37%
	SB Granular MoE	18.34%	26.78%	23.38%	22.83%
	PB Granular MoE	18.43%	27.05%	24.72%	23.40%
400B FineWeb → 200B German CC	Dense Baseline	19.28%	32.53%	23.99%	25.27%
	SB Switch MoE	21.76%	35.74%	25.21%	27.57%
	PB Switch MoE	20.73%	35.77%	23.01%	26.50%
	SB Granular MoE	22.61%	36.90%	<b>26.19%</b>	<b>28.57%</b>
	PB Granular MoE	22.70%	<b>37.23%</b>	23.01%	27.65%
400B FineWeb ∪ 200B German CC	Dense Baseline	20.05%	31.45%	24.85%	25.45%
	PB Granular MoE	21.33%	35.74%	25.70%	27.59%



## C.5 Training and Validation Loss

In the following sections, we present validation loss curves during and after pre-training and continual pre-training for all models in our study. Specifically, we report final validation loss in Table 5 and validation curves during training across figures 6, 7, and 8.

Table 5: **Final validation loss of MoEs and dense model after pre-training (English Web Data) and continual pre-training (Code & German Web Data).** As expected, we observe that all MoE transformers outperform the dense baseline during pre-training and continual pre-training with respect to validation loss. Moreover, we observe that MoEs forget marginally less than their dense counterparts. Together, these results show the continual learning abilities of MoEs are on par with dense models in terms of adaptation and are slightly superior in terms of forgetting, possibly due to their larger total parameter count.

Training Tokens	Model	Final Validation Loss				
		FineWeb	Stack	German	Forgetting	AVG
400B FineWeb (non-annealed)	Dense Baseline	2.881	4.028	3.741	–	–
	SB Switch MoE	2.711	3.861	3.495	–	–
	PB Switch MoE	2.699	3.872	3.451	–	–
	SB Granular MoE	2.664	3.690	3.404	–	–
	PB Granular MoE	2.653	3.715	3.370	–	–
400B FineWeb (annealed)	Dense Baseline	2.825	4.028	3.741	–	–
	SB Switch MoE	2.640	3.861	3.495	–	–
	PB Switch MoE	2.628	3.872	3.451	–	–
	SB Granular MoE	2.595	3.690	3.404	–	–
	PB Granular MoE	2.582	3.715	3.370	–	–
400B FineWeb → 200B Stack 30% Replay	Dense Baseline	2.939	1.026	–	0.059	1.982
	SB Switch MoE	2.757	0.944	–	0.046	1.850
	PB Switch MoE	2.749	0.945	–	0.050	1.847
	SB Granular MoE	2.708	0.925	–	0.044	1.816
	PB Granular MoE	2.699	0.924	–	0.046	1.811
400B FineWeb ∪ 200B Stack	Dense Baseline Union	2.866	1.050	–	–	1.958
	PB Granular MoE Union	2.630	0.935	–	–	1.782
400B FineWeb → 200B German 0% Replay	Dense Baseline	4.028	–	1.279	1.399	2.654
	SB Switch MoE	3.810	–	1.062	1.180	2.436
	PB Switch MoE	3.782	–	1.059	1.152	2.420
	SB Granular MoE	3.701	–	1.038	1.071	2.369
	PB Granular MoE	3.685	–	1.028	1.055	2.356
400B FineWeb → 200B German 40% Replay	Dense Baseline	2.946	–	1.367	0.066	2.157
	SB Switch MoE	2.749	–	1.142	0.039	1.946
	PB Switch MoE	2.741	–	1.129	0.042	1.935
	SB Granular MoE	2.701	–	1.118	0.037	1.910
	PB Granular MoE	2.690	–	1.099	0.037	1.895
400B FineWeb ∪ 200B German	Dense Baseline Union	2.938	–	1.390	–	2.164
	PB Granular MoE Union	2.669	–	1.120	–	1.895

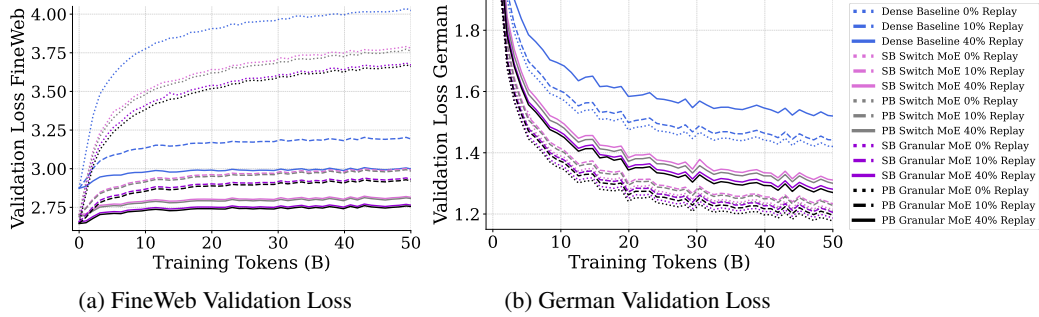


Figure 6: **Penalty-Balanced (PB) and Sinkhorn-Balanced (SB) Top- $k$  MoEs behave similarly to the FLOP-matched Dense baseline when being continually pre-trained with varying amounts of replay.** We continually pre-train MoEs and a dense baseline using varying amounts of replay: 0% (dotted curves), 10% (dashed curves), and 40% (full curves). We observe that replay substantially reduces forgetting for all models while slightly harming adaptation; that is, the effect of replay is the same for MoEs as for dense models.

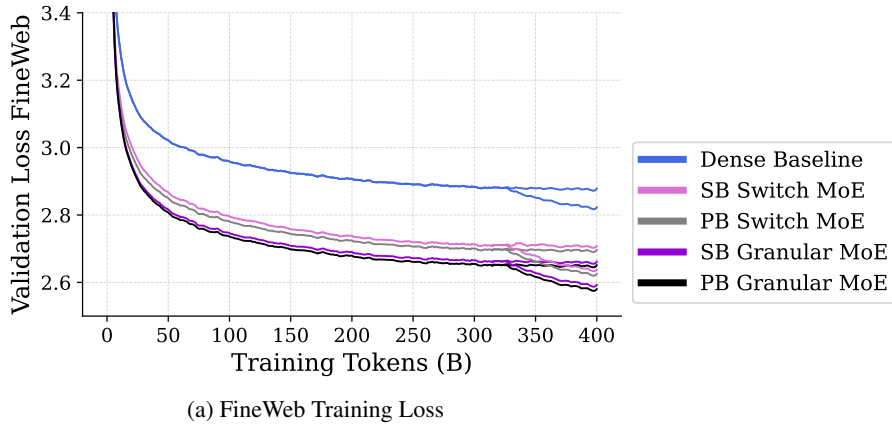


Figure 7: **Validation loss during initial pre-training on FineWeb with Infinite LR schedules.** We report decay and constant phases to completion. We observe that all MoE transformers stably decrease validation loss throughout pre-training, with MoEs improving over the dense model as expected. Interestingly, the PBT $k$  MoEs shows an incremental improvement over SBT $k$ .

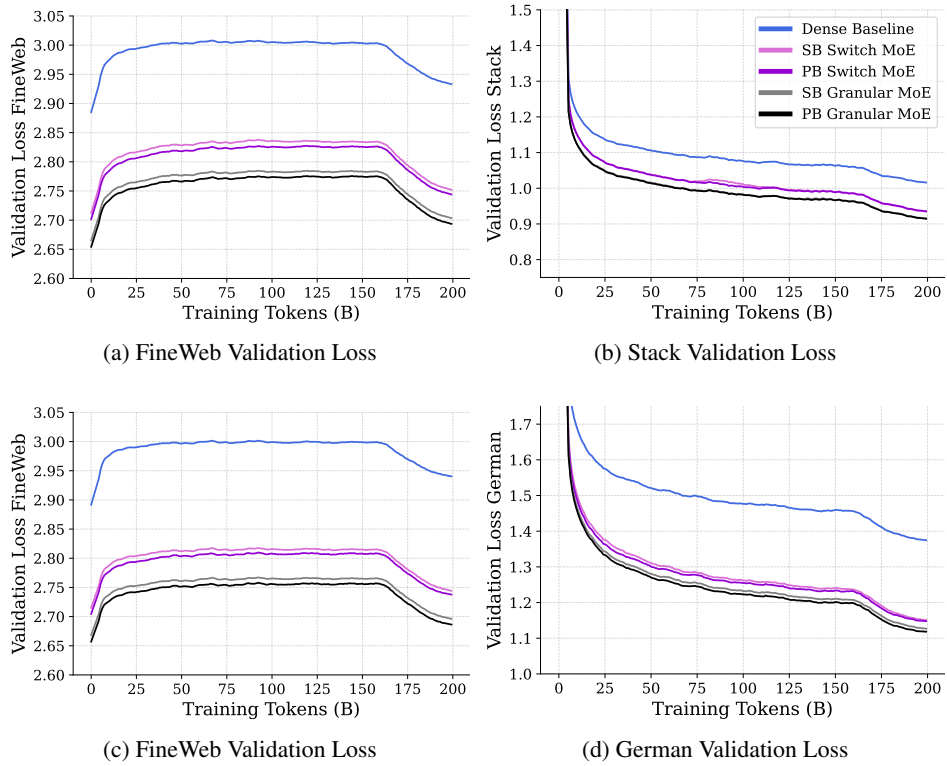


Figure 8: **Validation Loss on CPT and PT datasets during CPT.** Subfigures (a) and (c) report FineWeb validation loss, while subfigures (b) and (d) report Stack and German validation loss respectively for models trained on those datasets. We observe that all MoEs maintain their sample efficiency after the distribution shift, reaching a lower loss in many fewer iterations than the FLOP-matched dense baseline.

## C.6 Qualitative Analysis

In the following section, we present new metrics for analyzing the routing decisions of MoEs in our study and interpret how they change during continual pre-training. To accomplish this, we take checkpoints before and after continual pre-training and record their routing decisions, loss, routing imbalance, and a number of different metrics on 20M tokens of FineWeb test data (pre-training dataset), 20M tokens of German test data (continual pre-training dataset), and 20M tokens of Stack test data (continual pre-training dataset). Our results can be grouped into four main categories: 1) routing saturation analysis, 2) vocabulary specialization analysis, 3) expert co-activation analysis, and 4) routing imbalance analysis.

### C.6.1 Continual routing saturation analysis

We adapt the analysis of router saturation from [46] to the continual setting. Note that we will directly reproduce and slightly modify some lines from [46]’s definition of router saturation below for clarity and ease of passing from one paper’s notation to the other. Concretely, we define continual Continual Router Saturation as:

$$\text{Continual Router Saturation}(t, h, j) = \frac{1}{N} \sum_{i=1}^N \frac{|\mathcal{E}_i^{(\mathcal{T}_h)} \cap \mathcal{E}_i^{(\mathcal{T}_j)}|}{k}, \quad (2)$$

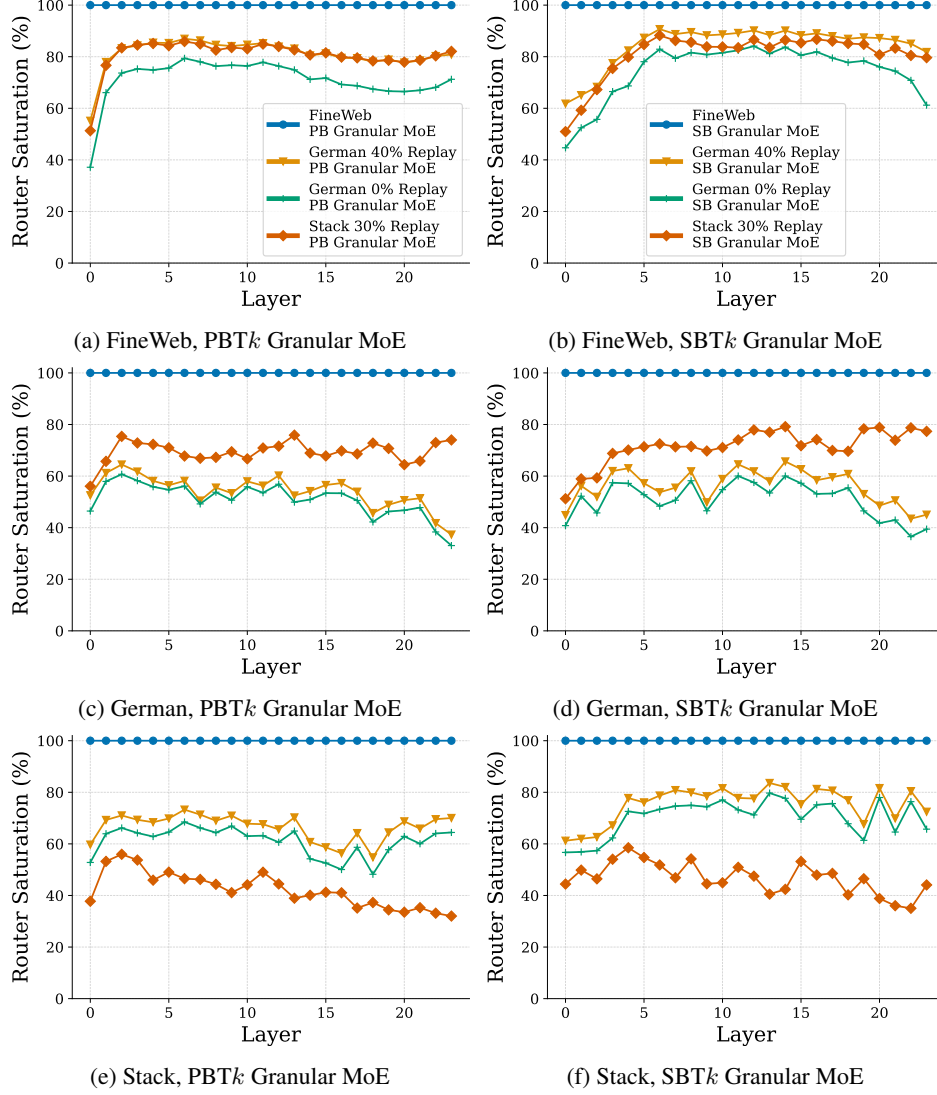
where:

- $\mathcal{T}_h$  and  $\mathcal{T}_j$ : The tasks being considered when selecting checkpoints. Note that  $h \leq j$ . In our case,  $j, h \in \{0, 1, 2\}$ , with  $\mathcal{T}_0$  designating the pre-training task (FineWeb) and  $\mathcal{T}_1, \mathcal{T}_2$  designating German and Stack continual pre-training tasks, respectively. While our experiments only consider one transition, in general, there may be many more.
- $N$ : The total number of tokens in the dataset.
- $k$ : The number of experts activated per input token.
- $\mathcal{E}_i^{(\mathcal{T}_h)}$ : The set of  $k$  experts activated for the  $i$ th token at the final checkpoint of the  $h$ th task.
- $\mathcal{E}_i^{(\mathcal{T}_j)}$ : The set of  $k$  experts activated for the  $i$ th token at the final checkpoint of the  $j$ th task.
- $|\mathcal{E}_i^{(\mathcal{T}_h)} \cap \mathcal{E}_i^{(\mathcal{T}_j)}|$ : The number of common experts activated for the  $i$ th token between the final checkpoints taken from the  $h$ th task and  $j$ th task.

Figures 9 and 10 consider  $h = 0$  and  $j \in \{0, 1\}$  for Granular (31 routed experts, 3 active, 1 shared) and Switch (8 routed experts, 1 active) MoEs respectively, thus comparing the checkpoint before continual pre-training with the checkpoint obtained afterward. The subfigures on the right report router saturation across model layers for PBT $k$  MoEs, while the subfigures on the left report the same for switch SBT $k$  MoEs. Each row reports router saturation on a different dataset. We make the following observations:

- (1) the first few layers are consistently among those with the lowest router saturation,
- (2) router saturation is lower for checkpoints trained on a given task  $h$  when it is measured with respect to tokens from  $h$ , and
- (3) the router saturation of models tested on their continual pre-training dataset seems to consistently decrease with a small slope as the layers index increases.

Observation (1) suggests that the early layers may undergo the most change during continual pre-training. Note that the trend of the first few layers having low router saturation is especially pronounced in subfigures (a) and (b), suggesting that most of the forgetting seen during continual pre-training may occur in the early layers. Observation (2) shows that MoEs change their routing decisions more to the distribution they are being trained on in the case of German and Stack, which is intuitive. Observation (3) suggests that layers closer to the final layer of the MoE must change more to adapt to the new distribution.



**Figure 9: Router saturation at the beginning of continual pre-training for Granular MoEs.** Subfigures (a,c,e) report layer-wise router saturation for PBTk MoEs, while subfigures (b,d,f) report router saturation for SBTk MoEs. (a) and (b) measure routing saturation with respect to test tokens from FineWeb, (c) and (d) measure router saturation with respect to test tokens from German, and (e) and (f) measure router saturation with respect to test tokens from Stack. We observe a few trends: 1) the first few layers are consistently among those with the lower router saturation, 2) router saturation is consistently lower for checkpoints CPT on the testing distribution showing that these checkpoints adapt more to that distribution, 3) the router saturation of models tested on their continual pre-training dataset seems to consistently decrease with a small slope as the layers index increases, and 4) the no-replay checkpoint consistently has lower router saturation than its 40% replay counterpart.



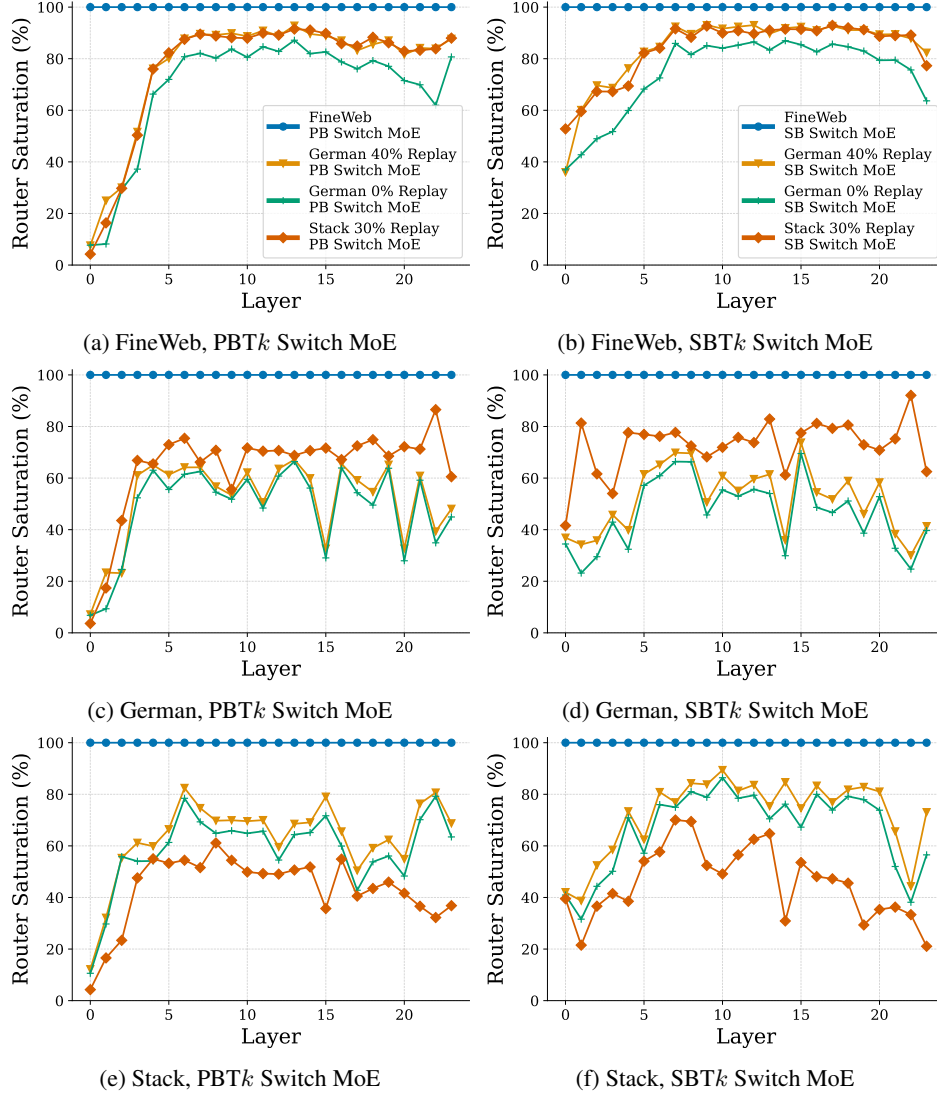


Figure 10: **Router saturation at the beginning of continual pre-training for Switch MoEs.** Subfigures (a,c,e) report layer-wise router saturation for PBT $k$  MoEs, while subfigures (b,d,f) report router saturation for SBT $k$  MoEs. (a) and (b) measure routing saturation with respect to test tokens from FineWeb, (c) and (d) measure router saturation with respect to test tokens from German, and (e) and (f) measure router saturation with respect to test tokens from Stack. We observe a few trends: 1) the first few layers are consistently among those with the lower router saturation, 2) router saturation is consistently lower for checkpoints CPT on the testing distribution showing that these checkpoints adapt more to that distribution, 3) the router saturation of models tested on their continual pre-training dataset seems to consistently decrease with a small slope as the layers index increases, and 4) the no-replay checkpoint consistently has lower router saturation than its 40% replay counterpart.

## C.6.2 Continual Vocabulary Specialization Analysis

We adapt the analysis of vocabulary specialization from [46] to the continual setting. Note that we will directly reproduce and slightly modify some lines from [46]’s definition of vocabulary specialization below for clarity and ease of passing from one paper’s notation to the other. Concretely, we define Vocabulary Specialization as:

$$\text{Vocabulary Specialization}(j, E_i, x) = \frac{N_{x, E_i, j}^{(k)}}{N_x}, \quad (3)$$

where:

- $E_i$ : The  $i$ th expert in an MoE layer.
- $j$ : A task index specifying which final checkpoint to use (e.g., specifying the final checkpoint after task 1, task 2,...).
- $x$ : The token ID being analyzed.
- $k$ : The number of experts considered (we use  $k=3$  for Granular MoEs and  $k=1$  for switch MoEs).
- $N_{j, x, E_i}^{(k)}$ : The number of times input data is routed to  $E_i$  for  $x$  when using the final checkpoint of task  $j$ .
- $N_{j, x}$ : The total number of times input data is routed across all experts for  $x$  and the final checkpoint of task  $j$ .

Vocabulary Specialization can, therefore, be calculated for each expert at every layer of the model and for each token in the model’s vocabulary. By assigning each token in the vocabulary to the expert that processes it the most frequently, we can then create a *one-to-many mapping* between experts and vocabulary entries for each layer of the MoE. Then, we can calculate the average vocabulary specialization of each expert by averaging over its assigned tokens and averaging across experts to measure specialization within a layer. To compare specialization across model checkpoints, we can re-use the *one-to-many mapping* of a previous checkpoint and measure how the specialization with respect to this mapping has changed during continual pre-training. Concretely, the continual vocabulary specialization (CVS) for an MoE layer  $l$  can be defined as follows:

$$\text{CVS}(j, h) = \frac{1}{N_E} \sum_{x \in \mathcal{V}} \text{Vocabulary Specialization}(h, E_{\alpha_{j, x}}, x) \quad (4)$$

$$\alpha_{j, x} := \arg \max_{i \in [N_E]} \{ \text{Vocabulary Specialization}(j, E_i, x) \} \quad (5)$$

- $N_E$ : The number of experts in an MoE layer  $l$ .
- $\mathcal{V}$ : The set of tokens in the model’s vocabulary (we use the Llama3 tokenizer).
- $h$ : A task index specifying the checkpoint from which to compute the mapping.
- $j$ : A task index specifying a final checkpoint that is used to compute the continual vocabulary specialization.

Note that the dataset of tokens used to compute the CVS is omitted for simplicity. However, the specialization of experts will depend on the distribution of the tokens because the same input token may be routed to different experts depending on the context within which it lives and the context will change depending on the distribution. For instance, the hidden representation of the word “for” in an English language corpus and a code corpus may differ wildly.

Figures 11 and 12 report the CVS of Granular and Switch MoEs, respectively. For each plot, the *one-to-many mapping*,  $\alpha_{j, x}$ , is created from the checkpoint pre-trained on FineWeb (e.g., the checkpoint we start continual pre-training from). All specializations are computed with respect to the input token. When evaluated on FineWeb, we observe across all architectures and balancing strategies that the first few layers for continually pre-trained models have lower continual vocabulary specialization than the pre-trained checkpoint, whereas subsequent layers have vocabulary specialization that closely matches that of the pre-trained checkpoint. This is even the case for the model that uses 0% replay,

1030 suggesting that MoEs learn routing policies during pre-training that are relatively unaffected by  
1031 continual pre-training. When evaluated on German and stack, we observe that all MoEs continually  
1032 pre-trained on those datasets have lower vocabulary specialization than models not trained on those  
1033 distributions, showing their adaptation. On German, the zero-replay model has the smallest CVS. We  
1034 hypothesize that this is the case because these models adapt the most to the German distribution and  
1035 happen to learn new routing patterns, distinct from the ones used on FineWeb. Contrasting the results  
1036 observed in subfigures (a) and (b) across Figures 11 and 12 to other subfigures, we observe that the  
1037 vocabulary specialization on the pre-training dataset only changes for the first few layers, while it  
1038 changes across all layers for the data seen during continual pre-training, even for the model that does  
1039 not utilize any replay. Contrasting this with the stronger performance of the no-replay model on  
1040 German and its poorer performance on FineWeb, the superior adaptation to German is correlated to  
1041 the change in vocabulary specialization throughout the model while the poorer performance on the  
1042 previous distribution is correlated with larger changes in vocabulary specialization in the first few  
1043 layers.

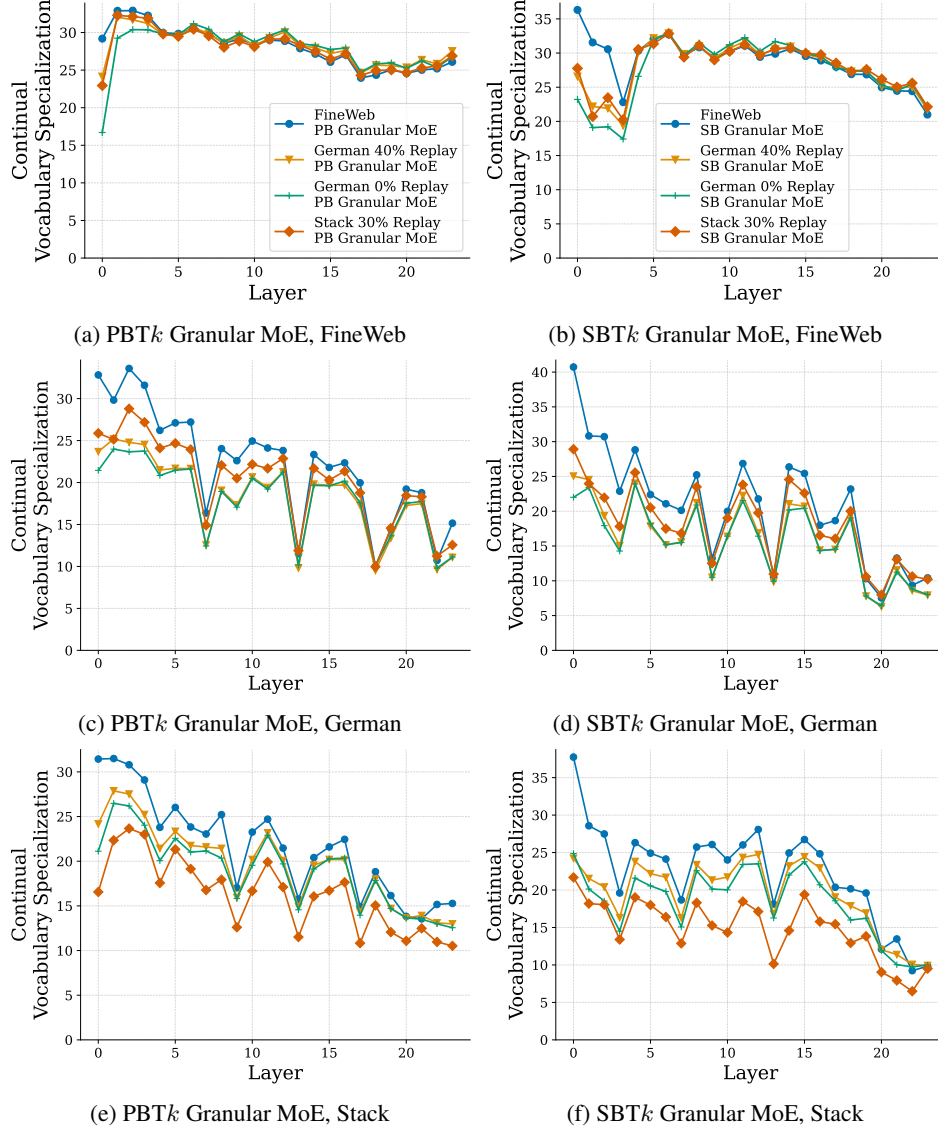


Figure 11: **Continual Vocabulary Specialization for Granular MoEs.** We report CVS for each MoE layer in the MoEs when testing models on FineWeb, German, and Stack. We observe that early layers deviate most from the checkpoint after pre-training, while later layers in the continually pre-trained MoEs nearly match the vocabulary specialization of their checkpoints after the first phase pre-training. This is even the case for the checkpoint that does not replay previous data, suggesting that vocabulary specialization for pre-training data is mostly determined during the initial pre-training phase.

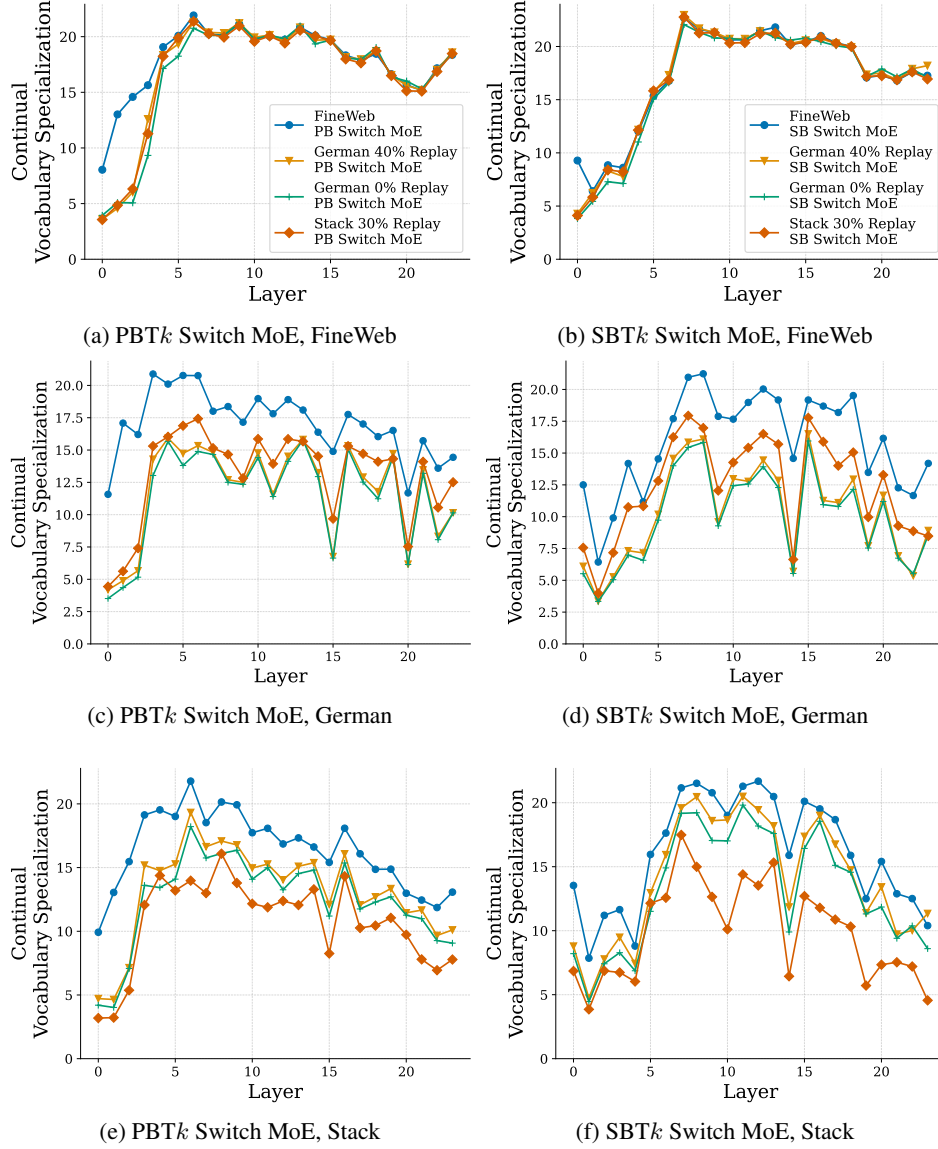


Figure 12: **Continual Vocabulary Specialization for Switch MoEs.** We report CVS for each MoE layer in the MoEs when testing models on FineWeb, German, and Stack. We observe that early layers deviate most from the checkpoint after pre-training, while later layers in the continually pre-trained MoEs nearly match the vocabulary specialization of their checkpoints after the first phase pre-training. This is even the case for the checkpoint that does not replay previous data, suggesting that vocabulary specialization for pre-training data is mostly determined during the initial pre-training phase.



### 1044 C.6.3 Continual expert co-activation analysis

1045 We adapt the analysis of expert co-activation from [46] to the continual setting. Note that we will  
 1046 directly reproduce and slightly modify some lines from [46]’s definition of expert co-activation below  
 1047 for clarity and ease of passing from one paper’s notation to the other. Concretely, we define Expert  
 1048 Co-activation as:

$$\text{Expert co-activation}(E_i, E_j) = \frac{N_{E_i, E_j}}{N_{E_i}}, \quad (6)$$

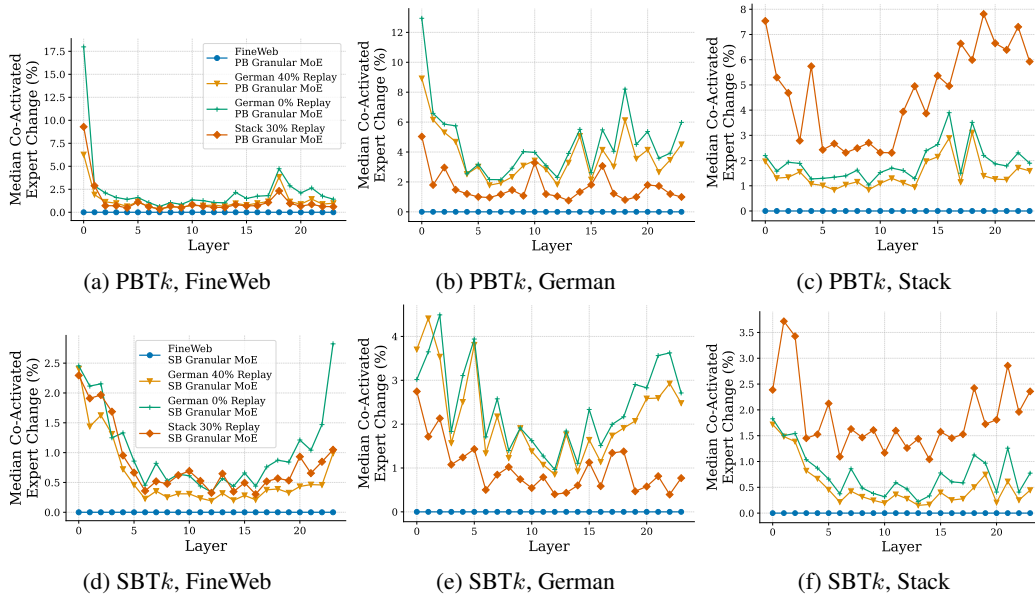
1049 where:

- 1050 •  $E_i$ : The first expert.
- 1051 •  $E_j$ : The second expert.
- 1052 •  $N_{E_i, E_j}$ : The number of times experts  $E_i$  and  $E_j$  are activated together.
- 1053 •  $N_{E_i}$ : The total number of times expert  $E_i$  is activated.

1054 The co-activation matrix  $C$  for any layer in the MoE can, therefore, be created by setting  $C_{i,j} =$   
 1055  $\text{Expert co-activation}(E_i, E_j)$ . Then, we can define the co-activation difference as follows:

$$\text{Co-activation Difference}(p, q) = |C^{(p)} - C^{(q)}|. \quad (7)$$

1056 Where  $|\cdot|$  is the coordinate-wise absolute value function. Each coordinate  $i, j$  of the co-activation  
 1057 difference measures the change in expert co-activation for experts  $i, j$  between final MoE checkpoints  
 1058 after tasks  $p$  and  $q$ , respectively. Taking statistics of the entries of the co-activation difference  
 1059 matrix allows us to measure how expert co-activation changes globally at each layer during continual  
 1060 pre-training.



**Figure 13: Layer-wise Median Router Co-activation Difference for Granular MoEs.** We report the median of the coordinates of the co-activation difference matrix between each model in the legend and its corresponding pre-trained checkpoint. We observe that on FineWeb, the Penalty-Balanced MoEs have the largest median differences overall and that they are most pronounced in the first two layers and layer 18. On German, we observe that the models continually pre-trained on German have the largest median differences and that the tendency for Penalty-Balanced MoEs to have large differences is maintained. On Stack, similar trends are observed.

1061 In Figure 13, we report the median of the coordinates of the co-activation difference matrix between  
 1062 each continually pre-trained Granular MoE in our study (the switch MoEs only activate a single expert

1063 so they have no co-activation) and its checkpoint after the initial pre-training phase. We observe  
1064 that when evaluated on the FineWeb test set, the Penalty-Balanced MoEs have the largest median  
1065 differences overall and that they are most pronounced in the first two layers and layer 18. When  
1066 evaluated on the German test set, we observe that the models continually pre-trained on German  
1067 have the largest median differences and that the tendency for Penalty-Balanced MoEs to have large  
1068 differences is maintained. When evaluated on the Stack test set, similar trends are observed.

1069 In Figures 14 and 15, we visualize a subset of the full expert co-activation matrices for Penalty-  
1070 Balanced and Sinkhorn-Balanced MoEs, respectively. Specifically, we show expert co-activations  
1071 for the 16 experts with the largest co-activation values. The left-most plots show the co-activation  
1072 matrix of the checkpoint continually pre-trained on FineWeb without decaying, the middle plots show  
1073 the co-activation matrix of the checkpoint after continually pre-training on Stack, and the rightmost  
1074 figures show the co-activation difference matrix. Subfigure (a) shows layer 0, (b) shows layer 11,  
1075 and (c) shows the final layer. We observe that the co-activation difference is the largest for layer  
1076 0 for both Penalty-Balanced and Sinkhorn-Balanced MoEs. Notably, for the Sinkhorn-Balanced  
1077 granular MoE’s pre-trained checkpoint, most of the co-activation weight is placed on the expert 15.  
1078 However, this strong weighting on expert 15 is attenuated during continual pre-training. In contrast,  
1079 the co-activations are more dispersed in the Penalty-Balanced MoE. For layers 11 and 23, there is  
1080 minimal change between pre-training and continual pre-training.

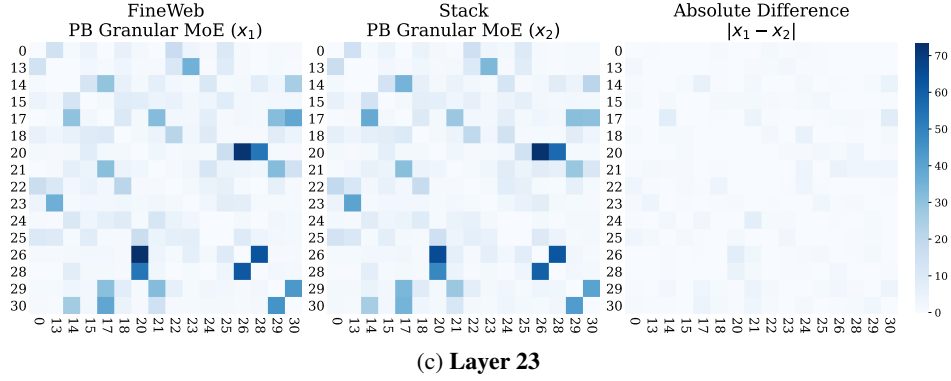
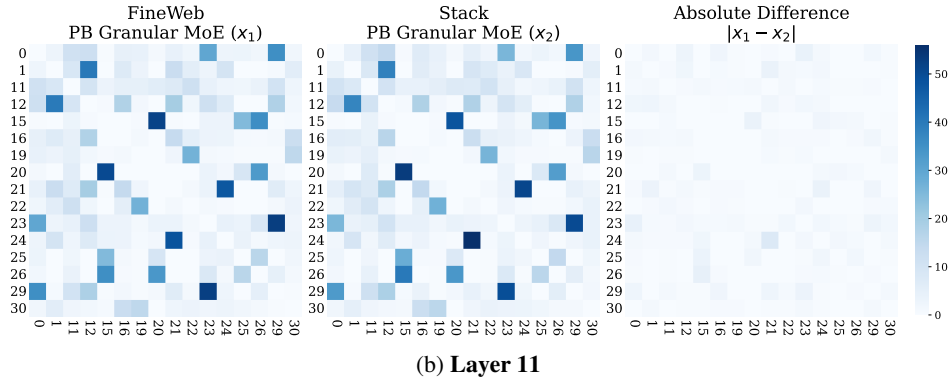
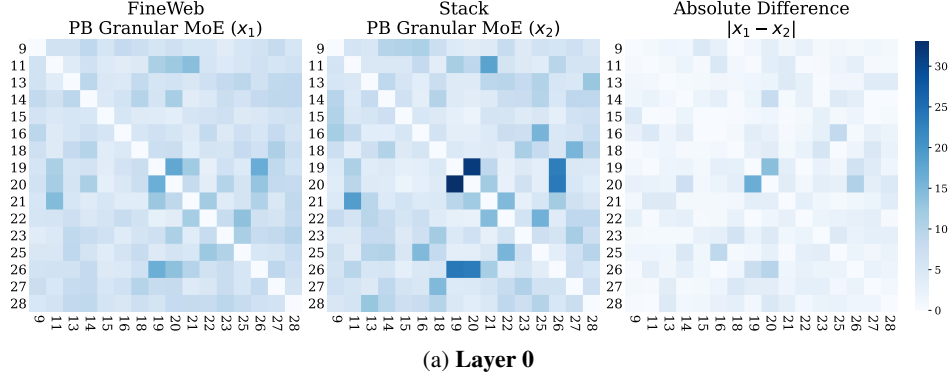
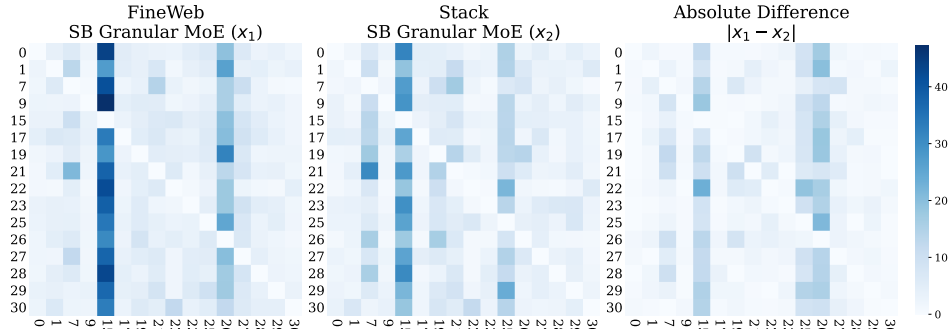
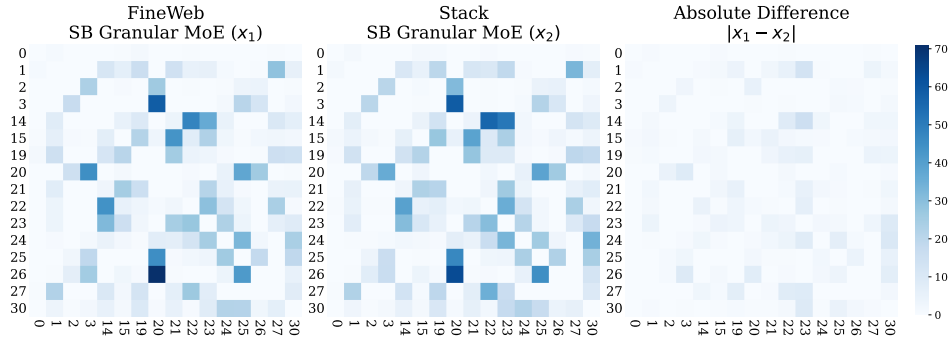


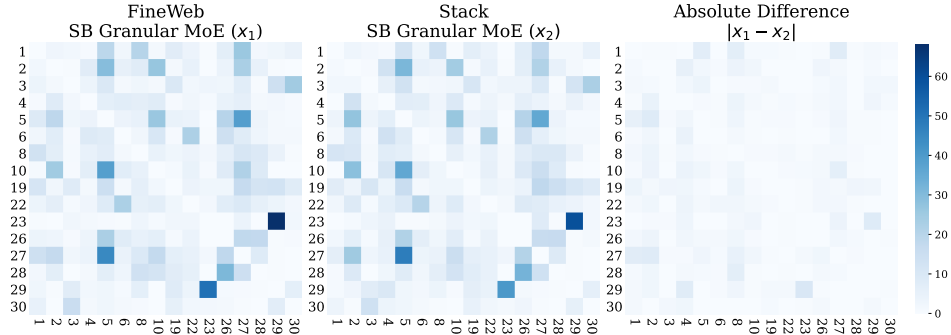
Figure 14: **FineWeb Router Co-activation Matrix for PB Granular MoE Continually Pre-trained on Stack.** The left-most plots show the Co-activation matrix of the checkpoint continually pre-trained on FineWeb without decaying, the middle plots show the Co-activation matrix of the checkpoint after continually pre-training on Stack, and the rightmost figures show the co-activation difference matrix. Subfigure (a) shows layer 0, (b) shows layer 11, and (c) shows the final layer. We observe that the co-activation difference is the largest layer 0, while the other layers change minimally.



(a) Layer 0



(b) Layer 11



(c) Layer 23

**Figure 15: FineWeb Router Co-activation Matrix for SB Granular MoE Continually Pre-trained on Stack.** The left-most plots show the Co-activation matrix of the checkpoint continually pre-trained on FineWeb without decaying, the middle plots show the Co-activation matrix of the checkpoint after continually pre-training on Stack, and the rightmost figures show the co-activation difference matrix. Subfigure (a) shows layer 0, (b) shows layer 11, and (c) shows the final layer. We observe that the co-activation difference is the largest layer 0, with most of the weight placed on expert 15. We observe that this strong weighting on expert 15 attenuated during continual pre-training. For layers, there is minimal change between pre-training and continual pre-training.

#### 1081 C.6.4 Continual routing imbalance analysis

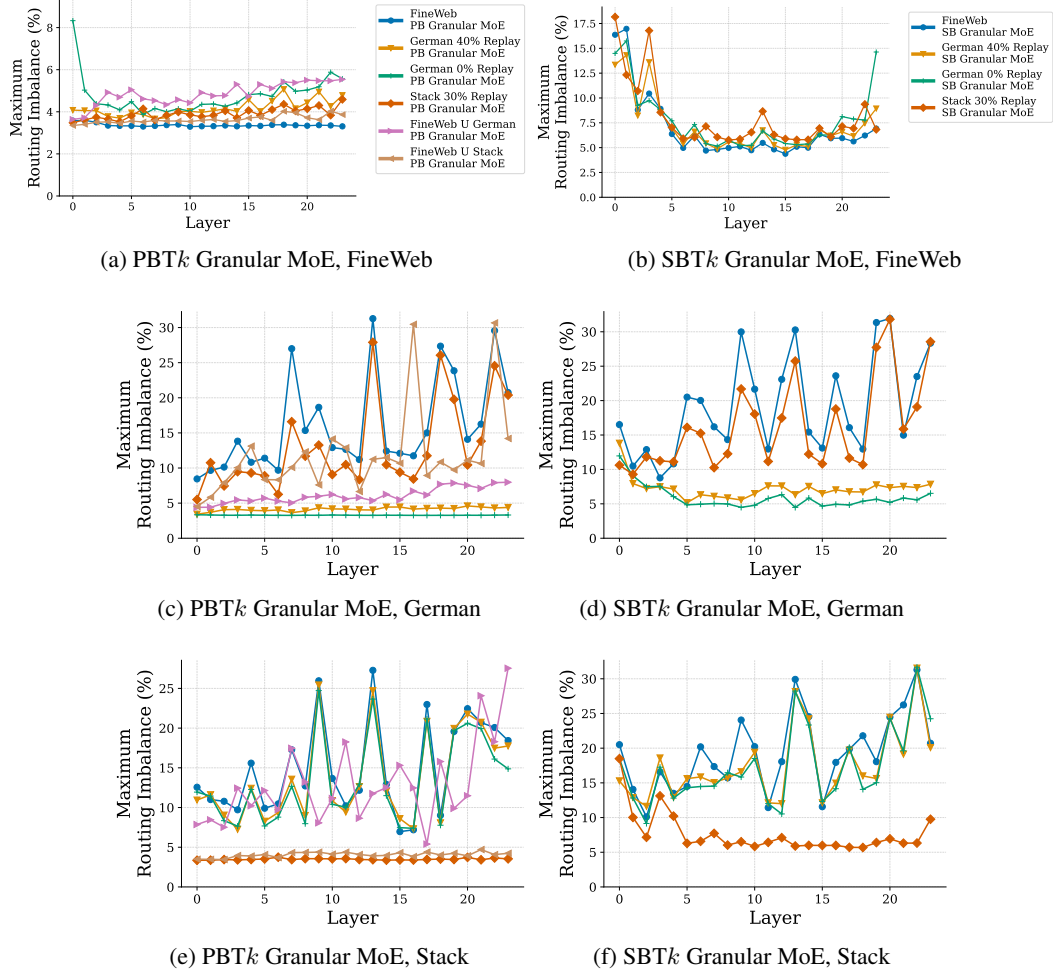
1082 While performance is one important axis of robustness to distribution shifts, maintaining a balanced  
 1083 load across experts is just as important for MoE foundation models. Without a balanced load, MoE  
 1084 transformers inferred using expert parallelism without token dropping (e.g., as is done for SOTA  
 1085 models [16, 69]) could be bottlenecked by the speed of a single accelerator that receives all the tokens,  
 1086 leading to underutilization of the hardware, lower throughput, and higher costs. To quantitatively  
 1087 assess the effect of distribution shift on load balance, we propose the maximum routing imbalance  
 1088 (MRI): the largest proportion of tokens routed to a single expert in a given MoE layer. Concretely,  
 1089 the MRI at a training iteration  $t$  and MoE layer  $j$  is defined as

$$\text{MRI}(t, j) := \max_{i \in [1, \dots, E]} \left[ \frac{\sum_{x \in B} \mathbb{1}\{i \in I_k(x)\}}{|B|} \right]. \quad (8)$$

1090 Where  $B$  is a set containing all tokens in a given batch,  $\mathbb{1}$  is the indicator function,  $E$  is the number of  
 1091 routed experts, and  $k$  is the number of active experts. *Since latency increases with computation, and,*  
 1092 *in an MoE layer, the computation required by a given device increases with the load of experts on*  
 1093 *that device, then MRI calculated with respect to routing decisions on a distribution is a proxy for the*  
 1094 *worst case latency of an MoE layer on the distribution.* We will use the MRI throughout the following  
 1095 sections to measure the effect of algorithmic changes to continual pre-training on routing imbalance.

1096 In Figures 16 and 17, we set  $t$  to be the final iteration of training for each model during pre-training  
 1097 and continual pre-training where it is applicable. The figures plot the layer identity on the x-axis and  
 1098 the MRI on the y-axis. The left column plots report the MRI of PBT $k$  MoEs, while the right column  
 1099 plots report MRI for SBT $k$  MoEs. Figures 16 shows Granular MoEs, while Figure 17 shows switch  
 1100 MoEs. For Granular MoEs, we observe that the MRI for Penalty-Balanced MoEs is consistently  
 1101 lower than for Sinkhorn-Balanced MoEs, that little increase in MRI on FineWeb is incurred during  
 1102 continual pre-training, even for the 0% replay model, and that MoEs become most unbalanced when  
 1103 seeing out-of-distribution data (e.g., see non-german models in (b) and non-code models in (c)). For  
 1104 Switch MoEs, we observe the MRI for Penalty-Balanced MoEs is similarly consistently lower than  
 1105 for Sinkhorn-Balanced MoEs, that similar to Granular MoEs little increase in MRI on FineWeb is  
 1106 incurred during continual pre-training, even for the 0% replay model, that switch MoEs become  
 1107 most unbalanced when seeing out-of-distribution data (e.g., see non-german models in (c,d) and  
 1108 non-code models in (e,f)), and that high MRI is prevalent in early layers independent of the training  
 1109 and testing distributions used, unlike for Granular MoEs. Contrasting these differences with the  
 1110 superior language modeling performance of granular MoEs, one could hypothesize that the unstable  
 1111 MRI observed in early layers for switch models that is not present in Granular MoEs may be a cause  
 1112 of the performance difference.





**Figure 16: Layer-wise Maximum Routing Imbalance (MRI) for Granular MoEs.** We report the MRI for each layer in the MoE as a percentage of all routing decisions made on a given dataset’s 20M token test set ((a,b)FineWeb, (c,d)German, and (e,f) Stack). The left column PBTk MoEs, while the left column reports results for SBTk MoEs. We observe that the MRI for Penalty-Balanced MoEs is consistently lower than for comparable Sinkhorn-Balanced MoEs, that little increase in MRI on FineWeb is incurred during continual pre-training, even for the 0% replay model (except for its first layer), and that MoEs become most unbalanced when seeing out-of-distribution data (e.g., see non-german models in (e,f) and non-code models in (c,d)).

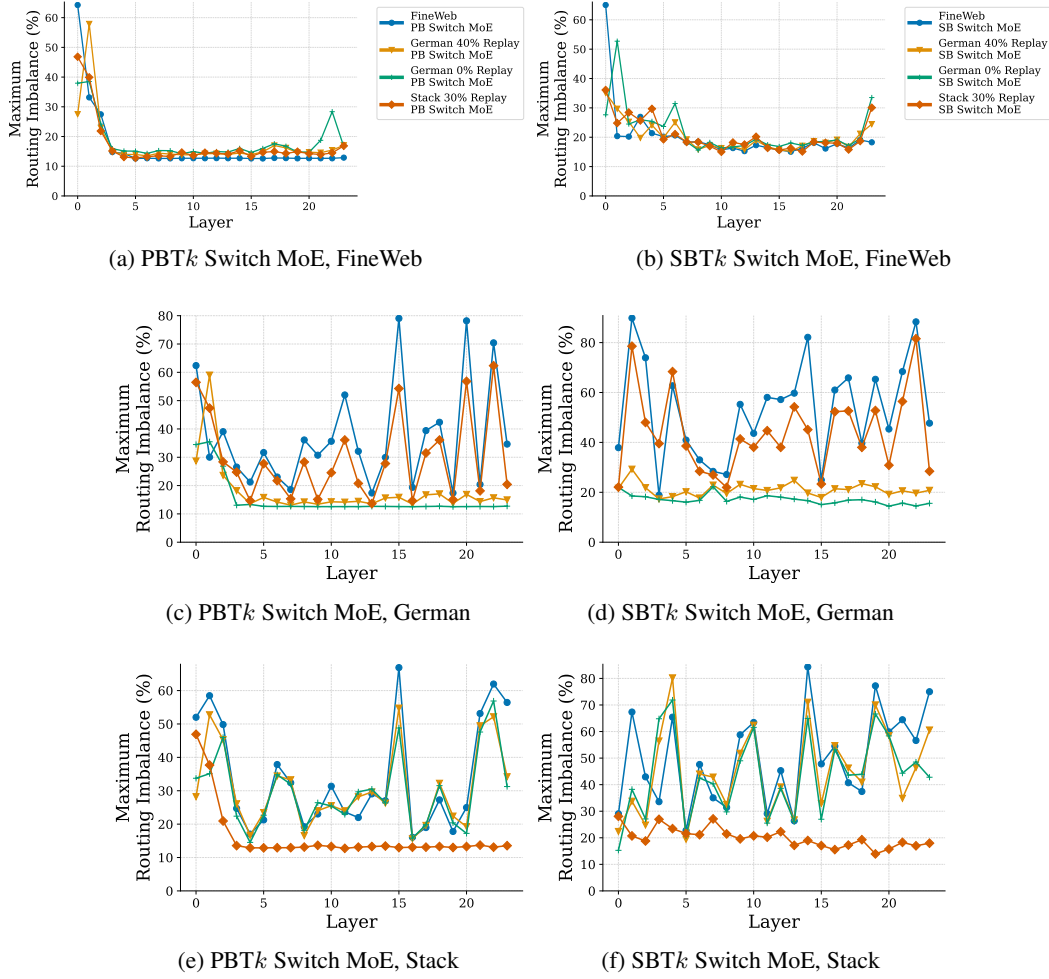


Figure 17: **Layer-wise Maximum Routing Imbalance (MRI) for Switch MoEs.** We report the MRI for each layer in the MoE as a percentage of all routing decisions made on a given dataset’s 20M token test set ((a,b) FineWeb, (c,d) German, and (e,f) Stack). The left column  $PBT_k$  MoEs, while the left column reports results for  $SBT_k$  MoEs. We observe that the MRI for Penalty-Balanced MoEs is consistently lower than for comparable Sinkhorn-Balanced MoEs, that little increase in MRI on FineWeb is incurred during continual pre-training, even for the 0% replay model (except for its first layer), that MoEs become most unbalanced when seeing out-of-distribution data (e.g., see non-german models in (e,f) and non-code models in (c,d)), and that high MRI is prevalent in early layers of Switch MoEs independent of the training and testing distributions used.

## 1113 C.7 Maximum routing imbalance of MoEs During Continual Pretraining.

1114 In the following section, we report on the maximum routing imbalance of MoEs during CPT.  
1115 Specifically, Figures 18 and 19 report routing immediately before and after the distribution shift,  
1116 while Figures 20 and 21 report MRI during training for Granular and switch MoEs.

1117 **Routing imbalance during training** By sparsely activating their weight matrices, MoEs experience  
1118 performance benefits over FLOP-matched dense models. However, this comes at the cost of increased  
1119 latency when the model’s forward pass is bottlenecked by the latency of a single expert. This can  
1120 become a problem if a router at any layer of the MoE chooses to dispatch a majority of the token  
1121 load to a particular expert. Therefore, we can estimate the impact of continual pre-training on MoE  
1122 latency by tracking the worst load imbalance at each layer of the MoE, which is the definition of the  
1123 maximum routing imbalance equation 1.

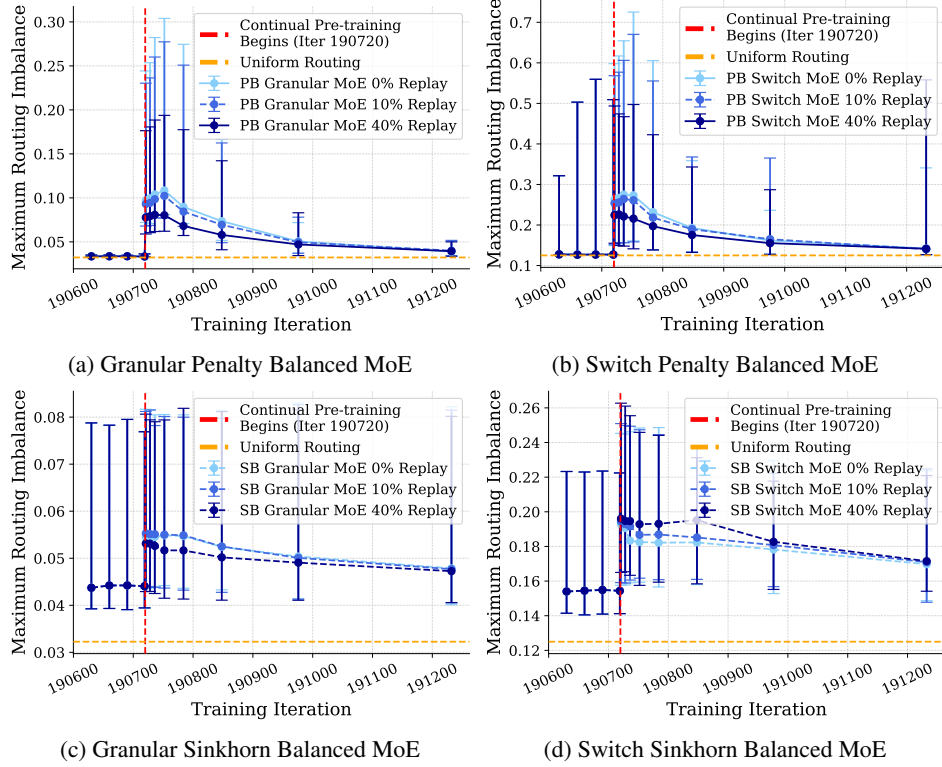
1124 In figures 20, and 21, we plot the MRI throughout pre-training (FineWeb) and continual pre-training  
1125 (German CC) for Switch and Granular MoEs, respectively. Subfigure (a) show the training time and  
1126 inference time MRI for PB MoEs, while subfigure (b) shows the training time MRI for SB MoEs and  
1127 subfigure (c) show the inference time MRI for SB MoEs. We distinguish between Sinkhorn Balanced  
1128 training and inference because the Sinkhorn Balancing algorithm is incompatible with autoregressive  
1129 generation, so in subfigure (c) we show MRI for SB models without the balancing step (e.g., what  
1130 would be used during autoregressive generation).

1131 For all MoEs, early layers (0-6) seem to have the largest MRI. For Switch and Granular MoEs  
1132 alike, we observe that SB routing follows a very similar pattern all throughout pre-training. During  
1133 the continual pre-training phase, we observe that this pattern changes slightly, actually becoming  
1134 more balanced throughout continual pre-training for both inference time and training time routing  
1135 imbalance. Turning our attention to the PB MoEs, we observe that Switch MoEs suffer from much  
1136 greater routing imbalance than their dense Granular counterparts in early layers. However, for most  
1137 layers of the PB switch MoE and all layers of the PB Granular MoE, the MRI quickly reaches a  
1138 smaller value than their SB counterparts during pre-training and continual pre-training showing that  
1139 PB MoEs are also robust to distribution shifts, but that the Granular MoE architecture is favorable for  
1140 continual pre-training

1141 *In summary, both PB and SB Top-k routing algorithms are robust to distribution shifts, with PB*  
1142 *initially being more perturbed by the distribution shifts, but recovering quickly to a better balance than*  
1143 *SB. These results demonstrate that using infinite LR schedules and replay is enough to continually*  
1144 *pre-train MoE LLMs without incurring a large increase in MRI.*

## 1145 C.8 Limitations

1146 While our study is comprehensive in scope and includes multiple state-of-the-art Mixture-of-Experts  
1147 (MoE) architectures and routing strategies, any empirical investigation operating under finite compute  
1148 constraints inevitably faces trade-offs. As such, although we evaluate four MoE configurations across  
1149 three large-scale datasets and two continual pre-training (CPT) transitions, our findings do not fully  
1150 capture the entire design space of MoE models. Specifically, we do not evaluate larger-scale MoEs  
1151 beyond the 2B parameter regime and we focus exclusively on a decoder-only transformer architecture  
1152 and do not study SSM-transformer hybrid variants, which may interact differently with routing  
1153 dynamics and distribution shifts. Moreover, while we study two of the most widely adopted Top- $k$   
1154 routing algorithms—Penalty-Balanced and Sinkhorn-Balanced routing—we do not exhaustively  
1155 evaluate alternative mechanisms such as hash routing [55]. Different routing approaches may exhibit  
1156 different trade-offs with respect to performance, stability, and generalization under CPT.



**Figure 18: Sinkhorn-Balanced (SB) and Penalty-balanced (PB) Top- $k$  MoEs show little change in training-time maximum routing imbalance as a result of adjusting the replay percentage.** We report the median MRI observed across MoE layers with min and max error bars shortly before and following the distribution shift when continually pre-training the MoEs on German CC. We observe that independent of the replay percentage used, the MoEs recover pre-training level median MRI within 1000 iterations of continual pre-training. However, replay does mitigate the increase in MRI caused by the distribution shift to a small extent in PB MoEs. In contrast, the SB MoEs are quite robust to the distribution shift and their routing patterns seem to be invariant to the replay percentage used.

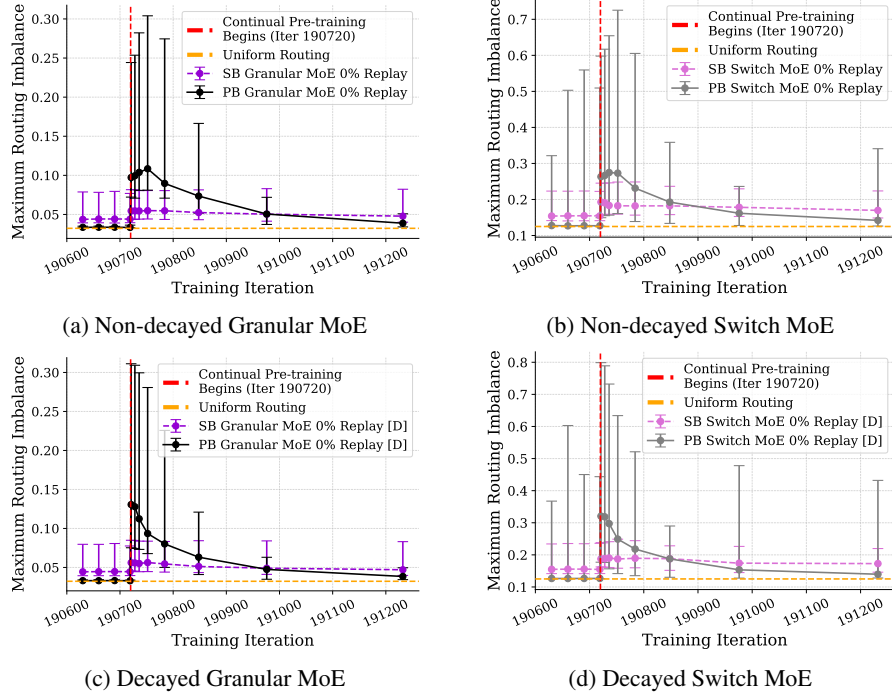
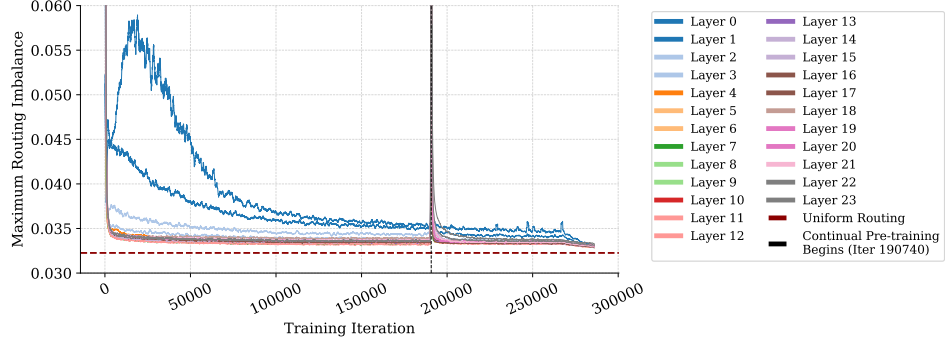
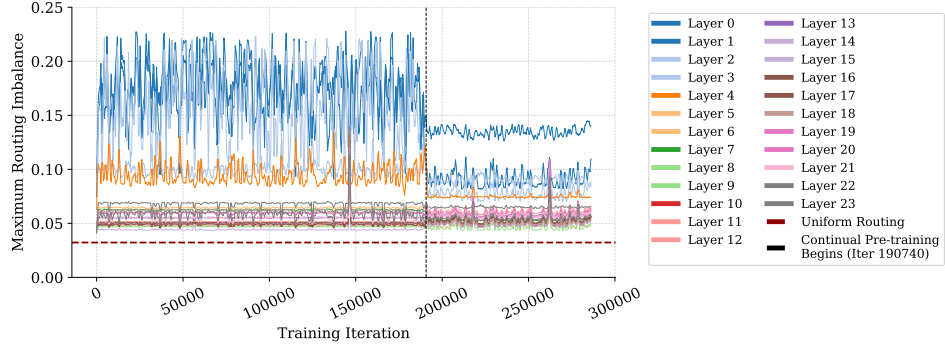


Figure 19: **Decayed Penalty-Balanced (PB) Top- $k$  MoEs have slightly higher MRI during distributions shifts than their non-decayed counterparts.** We report the median MRI observed across MoE layers with min and max error bars shortly before and following the distribution shift when continually pre-training the MoEs on German CC. We observe that all SB MoE keep a stable MRI throughout the distribution shift, showing that they are mostly unaffected. In contrast, the PB checkpoints suffer from strong routing imbalance after the distribution shift but recover quickly, with the decayed checkpoints reaching a marginally higher MRI.

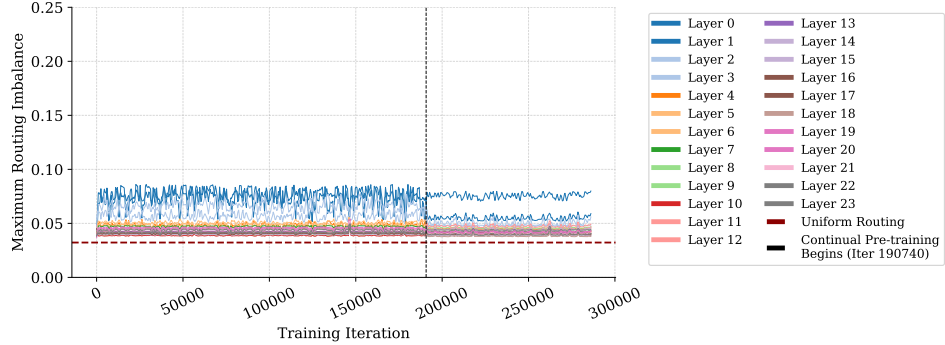




(a) Inference time & train time MRI, PB Granular MoE

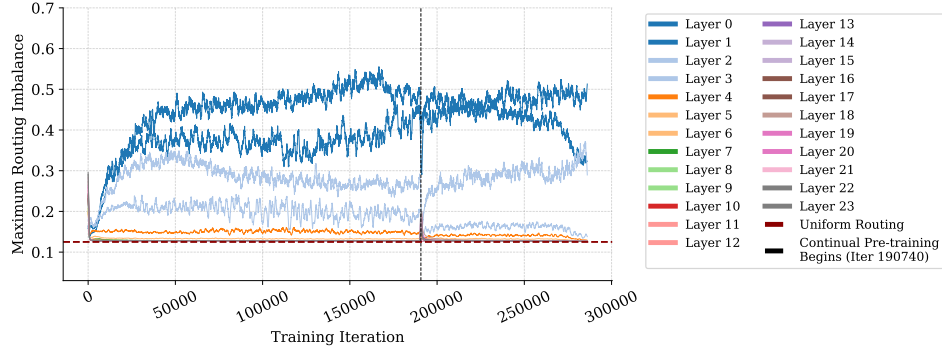


(b) Inference time MRI, SB Granular MoE

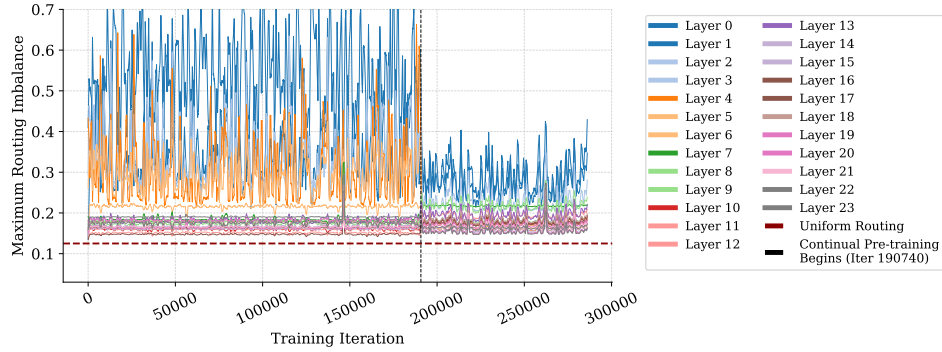


(c) Train time MRI, SB Granular MoE

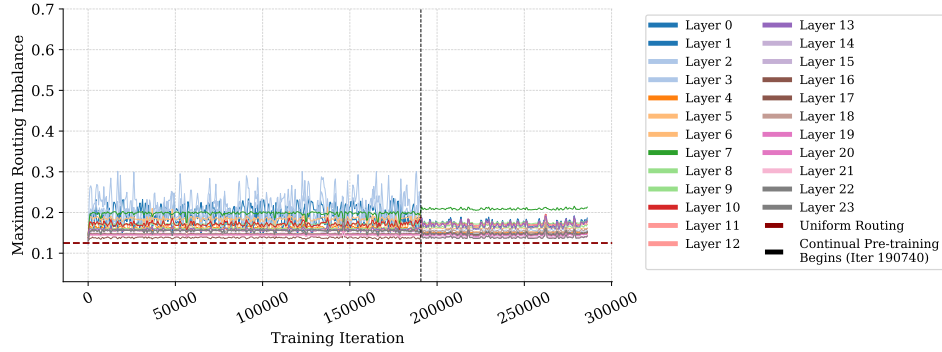
**Figure 20: Training time and inference time MRI for Granular MoEs throughout pre-training and continual pre-training.** We show layer-wise maximum routing imbalance during pre-training and continual pre-training. While Penalty-Balanced MoEs have the same routing dynamics at training and inference time, Sinkhorn balancing is incompatible with autoregressive generation, so we show both inference-time and train-time MRI for SB models. We observe that early layers in the MoE consistently have the largest MRI for both PB and SB MoEs, the MRI of PB MoEs is much better behaved, and after the distribution shift, the MRI of SB models becomes more stable.



(a) inference time & train time MRI, PB Switch MoE



(b) inference time MRI, SB Switch MoE



(c) Train time MRI, SB Switch MoE

**Figure 21: Training time and inference time MRI for Switch MoEs throughout pre-training and continual pre-training.** We show layer-wise maximum routing imbalance during pre-training and continual pre-training. While Penalty-Balanced MoEs have the same routing dynamics at training and inference time, Sinkhorn balancing is incompatible with autoregressive generation, so we show both inference-time and train-time MRI for SB models. We observe that early layers in the MoE consistently have the largest MRI for both PB and SB MoEs, the MRI of PB MoEs is much better behaved, and after the distribution shift, the MRI of SB models becomes more stable.

## 1157 D Dataset Sizes and Sampling Proportions

1158 In the following section, we report the training tokens used, training dataset sizes, and sampling  
 1159 proportions. Specifically, Table 6 reports the amount of data and its exact composition used for  
 1160 different pre-training phases. Tables 7, 8, and 9 report the amount of training tokens and sampling  
 1161 proportions used for FineWeb, German, and Stack, respectively.

Table 6: **Pre-training and Continual Pre-training Tokens.** We report the training tokens for all different model training configurations in this paper. During continual pre-training, each batch contains a proportion of replay tokens from the pre-training dataset and new tokens from the continual pre-training dataset.

Phase	Training Tokens	New Tokens	Replay Tokens
Pre-training	400B FineWeb	400B	–
Continual Pre-training	400B FineWeb → 200B Stack 30% Replay	140B	60B
	400B FineWeb → 200B German 40% Replay	120B	80B
	400B FineWeb → 200B German 0% Replay	200B	–

Table 7: **FineWeb CC: Train, Val, and Test dataset sizes used in our experiments.** For the purposes of our study, we create a more manageable subset of FineWeb by subsampling each Common Crawl dump within FineWeb into smaller subsets. We then sample proportional to the sizes of each subset. We report the full size of the subset we sample from during training (note, we only train on 400B tokens of this subset). The exact sizes and sampling proportions of each split are committed as they span more than a page, but can be made available upon request.

Source	Train Tokens (B)	Test Tokens (B)	Val. Tokens (B)	Sampling Weight
FineWeb CC	2916.650	26.442	26.426	1.000

Table 8: **German CC: Train, Val, and Test dataset sizes used in our experiments.**

Source	Train Tokens (B)	Test Tokens (B)	Val. Tokens (B)	Sampling Weight
German CC	169.291	0.489	0.491	1.000

Table 9: **Stack: Train, Val, and Test dataset sizes used in our experiments.**

Source	Training Tokens (B)	Test Tokens (B)	Val. Tokens (B)	Sampling Weights
YAML	9.039	0.613	0.609	0.017
Java	19.730	0.587	0.587	0.174
C	17.988	0.594	0.597	0.159
Markdown	21.699	0.477	0.474	0.017
PHP	16.660	0.450	0.447	0.146
C#	9.245	0.552	0.553	0.084
JSON	120.669	0.709	0.695	0.017
TypeScript	6.892	0.418	0.414	0.063
C++	13.998	0.538	0.539	0.124
Python	15.898	0.458	0.457	0.200
<b>Total</b>	251.819	5.396	5.372	1.000

## 1162 E Model Hyperparameters

1163 The following section outlines the hyperparameters used to train the MoEs and dense transformers in  
 1164 our study. Specifically, Table 10 reports the hyperparameters of the schedules and Table 11 reports  
 1165 the model hyperparameters. We also show an example infinite learning rate schedule in Figure 22.

Table 10: **Hyperparameters of LR schedules.** All models used the same LR schedule hyperparameters. We refer the readers to [26] section 7.2 for a more thorough explanation of these schedules.

Description	Value
<b>Pre-training</b>	
Schedule Type	<i>CosineInf</i>
Total Iterations	192720
Max learning rate ( $\eta_{max}$ )	$3 \cdot 10^{-4}$
Min learning rate ( $\eta_{min}$ )	$3 \cdot 10^{-5}$
Constant learning rate ( $\eta_{const}$ )	$1.65 \cdot 10^{-4}$
Warmup percent ( $T_{warmup}$ )	1
Cooldown iters percent ( $T_{cd}$ )	70
Constant iters percent ( $T_{ann}$ )	0.10
<b>Continual Pre-training</b>	
Schedule Type	<i>CosineInf</i>
Total Iterations	95370
Max learning rate ( $\eta_{max}$ )	$3 \cdot 10^{-4}$
Min learning rate ( $\eta_{min}$ )	$3 \cdot 10^{-5}$
Constant learning rate ( $\eta_{const}$ )	$1.65 \cdot 10^{-4}$
Warmup percent ( $T_{warmup}$ )	1
Cooldown iters percent ( $T_{cd}$ )	0
Constant iters percent ( $T_{ann}$ )	80
<b>Full re-training</b>	
Schedule Type	<i>Cosine Annealing</i>
Total Iterations	288090
Max learning rate ( $\eta_{max}$ )	$3 \cdot 10^{-4}$
Min learning rate ( $\eta_{min}$ )	$3 \cdot 10^{-5}$
Warmup percent ( $T_{warmup}$ )	1
<b>Continual Pre-training Ablation (Section C.1)</b>	
Schedule Type	<i>Cosine Annealing</i>
Total Iterations	95370
Max learning rate ( $\eta_{max}$ )	$3 \cdot 10^{-4}$
Min learning rate ( $\eta_{min}$ )	$3 \cdot 10^{-5}$
Warmup percent ( $T_{warmup}$ )	1

Table 11: **Hyperparameters of our Moes and Dense Transformer.**

Description	Value
<b>MoE Transformers Common</b>	
Active Parameters	571, 148, 288
Parameters	2, 025, 236, 480
Non-Embedding Parameters	1, 893, 902, 336
<b>MoE SM-FFN</b>	
Shared Experts	1
Active Experts	3
Routed Experts	31
Total Experts	32
FFN Intermediate Size	704
<b>MoE R-FFN</b>	
Shared Experts	0
Active Experts	1
Routed Experts	8
Total Experts	8
FFN Intermediate Size	2816
<b>Top-k</b>	
Z-loss Coeff.	0.001
AUX-loss Coeff.	0.01
<b>Sinkhorn</b>	
Tolerance	0.01
<b>Dense Transformer</b>	
Parameters	571, 148, 288
Non-Embedding Parameters	439, 814, 144
Num attention heads	16
<b>Common</b>	
Num layers	24
Hidden size	1024
FFN Hidden size	2816
FFN Type	GeGLU
Optimizer	AdamW
$\beta_1, \beta_2$	0.9, 0.95
Batch size	1024
Sequence length	2048
Hidden activation	GeLU
Weight decay	0.1
Gradient clipping	1.0
Decay	Cosine
Positional embedding	Rotary
GPT-J-Residual	True
Weight tying	False
Vocab Size	128000
Rotary PCT	0.25

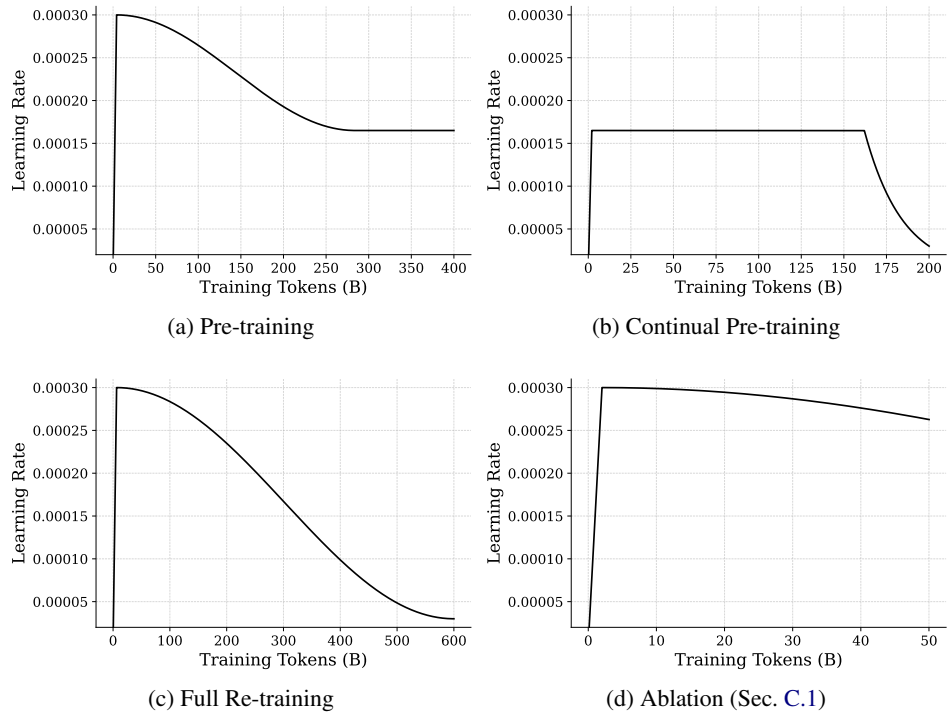


Figure 22: **Illustrated Learning rate schedules used for (a) pre-training, (b) continual pre-training, (c) full re-training, and the (d) rewarming ablation of Sec. C.1.** The exact hyperparameters of these schedules are reported in table 10.

## 1166 NeurIPS Paper Checklist

### 1167 1. Claims

1168 Question: Do the main claims made in the abstract and introduction accurately reflect the  
1169 paper's contributions and scope?

1170 Answer: [\[Yes\]](#)

1171 Justification: The abstract and introduction clearly articulate the paper's goals: evaluating the  
1172 robustness of MoE routers in CPT settings. The study thoroughly investigates performance  
1173 and load-balancing behavior across multiple MoE variants under realistic CPT scenarios,  
1174 aligning well with the claimed contributions.

1175 Guidelines:

- 1176 • The answer NA means that the abstract and introduction do not include the claims  
1177 made in the paper.
- 1178 • The abstract and/or introduction should clearly state the claims made, including the  
1179 contributions made in the paper and important assumptions and limitations. A No or  
1180 NA answer to this question will not be perceived well by the reviewers.
- 1181 • The claims made should match theoretical and experimental results, and reflect how  
1182 much the results can be expected to generalize to other settings.
- 1183 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
1184 are not attained by the paper.

### 1185 2. Limitations

1186 Question: Does the paper discuss the limitations of the work performed by the authors?

1187 Answer: [\[Yes\]](#)

1188 Justification: Section acknowledges the limitations of our study.

1189 Guidelines:

- 1190 • The answer NA means that the paper has no limitation while the answer No means that  
1191 the paper has limitations, but those are not discussed in the paper.
- 1192 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1193 • The paper should point out any strong assumptions and how robust the results are to  
1194 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
1195 model well-specification, asymptotic approximations only holding locally). The authors  
1196 should reflect on how these assumptions might be violated in practice and what the  
1197 implications would be.
- 1198 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
1199 only tested on a few datasets or with a few runs. In general, empirical results often  
1200 depend on implicit assumptions, which should be articulated.
- 1201 • The authors should reflect on the factors that influence the performance of the approach.  
1202 For example, a facial recognition algorithm may perform poorly when image resolution  
1203 is low or images are taken in low lighting. Or a speech-to-text system might not be  
1204 used reliably to provide closed captions for online lectures because it fails to handle  
1205 technical jargon.
- 1206 • The authors should discuss the computational efficiency of the proposed algorithms  
1207 and how they scale with dataset size.
- 1208 • If applicable, the authors should discuss possible limitations of their approach to  
1209 address problems of privacy and fairness.
- 1210 • While the authors might fear that complete honesty about limitations might be used by  
1211 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
1212 limitations that aren't acknowledged in the paper. The authors should use their best  
1213 judgment and recognize that individual actions in favor of transparency play an impor-  
1214 tant role in developing norms that preserve the integrity of the community. Reviewers  
1215 will be specifically instructed to not penalize honesty concerning limitations.

### 1216 3. Theory assumptions and proofs

1217 Question: For each theoretical result, does the paper provide the full set of assumptions and  
1218 a complete (and correct) proof?



Answer: [NA]

Justification: The paper is entirely empirical and does not include theoretical results requiring proofs.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The paper includes detailed training schedules, model architectures, datasets used, and evaluation protocols. Moreover, detailed Hyperparameter settings are included in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: At the time of writing, we are unable to provide access to the code due to an ongoing internal review process

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Section A.3 and Appendix E contain detailed model and optimizer configurations, replay ratios, and training schedules. These details enable replication.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: The paper does not report error bars or variance across multiple runs. Results are reported as point estimates due to high computational cost.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: The paper mentions model scale (2B parameters), token count (600B), and that experiments were conducted using A100 GPU clusters. Detailed compute resource estimates are included in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The paper follows NeurIPS ethical guidelines, uses public data, and does not present privacy, fairness, or safety risks.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: Our research furthers the field of machine learning, but does not directly have a societal impact.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper does not release pretrained models and does not pose any significant risk of misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The paper cites all datasets and prior models used.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets? item[] Answer: [NA]

Justification: The paper does not introduce new datasets, models, or other assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve human subjects or crowdsourcing.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The study does not involve human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- 1480
- 1481
- 1482
- 1483
- 1484
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
  - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

1485

## 16. Declaration of LLM usage

1486

1487

1488

1489

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

1490

Answer: [NA]

1491

1492

Justification: LLMs are the subject of the paper, not a tool used in producing it. The paper was not written using an LLM in any substantive research capacity.

1493

Guidelines:

- 1494
- 1495
- 1496
- 1497
- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
  - Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.