

MoFE: Mixture of Frozen Experts Architecture

Jean Seo, Jaeyoon Kim, Hyopil Shin

Seoul National University

{seemdog, toscour345, hpshin}@snu.ac.kr

Abstract

We propose the Mixture of Frozen Experts (MoFE) architecture, which integrates Parameter-efficient Fine-tuning (PEFT) and the Mixture of Experts (MoE) architecture to enhance both training efficiency and model scalability. By freezing the Feed Forward Network (FFN) layers within the MoE framework, MoFE significantly reduces the number of trainable parameters, improving training efficiency while still allowing for effective knowledge transfer from the expert models. This facilitates the creation of models proficient in multiple domains. We conduct experiments to evaluate the trade-offs between performance and efficiency, compare MoFE with other PEFT methodologies, assess the impact of domain expertise in the constituent models, and determine the optimal training strategy. The results show that, although there may be some trade-offs in performance, the efficiency gains are substantial, making MoFE a reasonable solution for real-world, resource-constrained environments.

1 Introduction

Large Language Models (LLMs) showcase significant advancements in natural language understanding and generation. LLMs are characterized by their immense size, often consisting of at least one billion parameters. The substantial size of LLMs is understandable given the scaling law suggested by Kaplan et al. (2020), which indicates that performance on the cross-entropy loss improves predictably with increased model size, data, and computational power. However, their immense size poses a resource challenge, requiring substantial computational memory and vast amounts of data, making development and deployment difficult to afford.

To address this, developing efficient LLMs that maintain high performance has become crucial. Efforts include **(1) Efficient Training Methodologies** like Parameter-efficient Fine-tuning (PEFT)

and **(2) Efficient Model Scaling Methodologies** such as the Mixture of Experts (MoE) architecture.

In this research, we propose the Mixture of Frozen Experts (MoFE) architecture, combining both approaches for a more efficient and affordable model. MoFE leverages MoE’s benefits while reducing computational requirements through freezing the FFN blocks. Our experiments demonstrate that, despite a trade-off between performance and efficiency compared to full fine-tuning, MoFE outperforms other PEFT methods, requiring the least training time while achieving the highest performance. Additionally, MoFE shows effective knowledge transfer from its constituent models, highlighting the potential for using pre-existing domain expertise models with minimal further training.

2 Related Work

Primary strategies for efficient model training are PEFT and quantization. PEFT includes techniques like prompt-tuning (Lester et al., 2021), adapters (Houlsby et al., 2019; Tomanek et al., 2021), LoRA (Hu et al., 2021), and DoRA (Liu et al., 2024), all designed to reduce computational demands. Quantization (Jacob et al., 2017) maps model weights to lower-precision formats for efficiency, and Dettmers et al. (2023) introduced QLoRA, combining LoRA with quantization.

The Mixture of Experts (MoE) architecture (Fedus et al., 2022; Shazeer et al., 2017; Komatsuzaki et al., 2023) is another efficient scaling method that gained attention with Mixtral 8X7B (Jiang et al., 2024), which integrates eight Mistral 7B models (Jiang et al., 2023) and outperforms Llama-2 70B (Touvron et al., 2023) despite being smaller. Following Mixtral 8X7B, other MoE-based models, including OpenMoE (Xue et al., 2024), Jamba (Lieber et al., 2024), BiMediX (Pieri et al., 2024), and BioMistral (Labrak et al., 2024), have been developed.

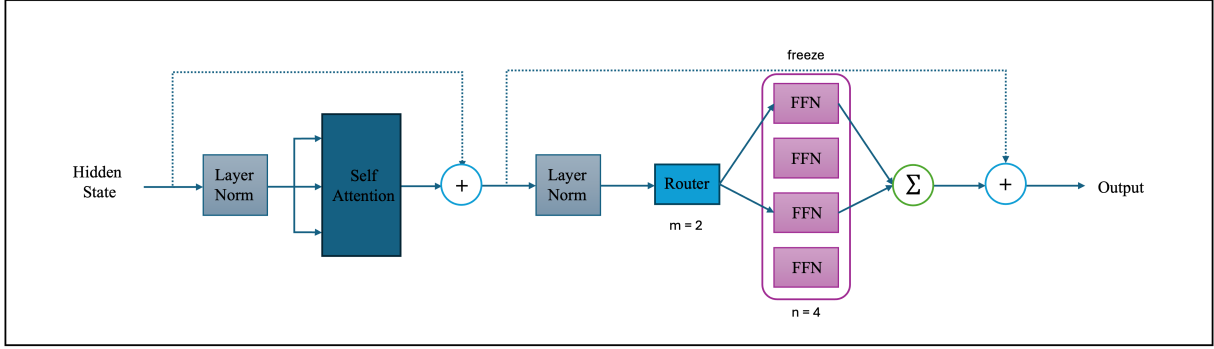


Figure 1: Mixture of Frozen Experts Architecture. In this example figure, the router uses 2 Feed Forward Network (FFN) blocks at each time step ($m = 2$), and there are 4 FFN blocks, or expert models, used ($n = 4$). In MoFE, the FFN blocks are frozen, so only the remaining parameters are updated. This makes the training process significantly more lightweight, regardless of the number of expert models integrated into the architecture.

3 MoFE

3.1 Architecture

We create a MoE model through the Mixtral architecture using mergekit (Goddard et al., 2024). The Mixtral architecture includes three components: the base model, the expert model, and the router. Here, the expert model provides the Feed Forward Network (FFN) layers, while the base model supplies other components like self-attention layers. In our experiments in Section 4, the models used as the base and expert models all have TinyLlama (Zhang et al., 2024), a pretrained model with 1.1 billion parameters, as the foundational model. As shown in Figure 1, the router (or gate) determines the number of FFN blocks used per time step, set to 2 ($m = 2$) in all experiments. In the proposed MoFE architecture, FFN blocks are frozen, while only the router and other parts are updated, keeping the trainable parameter size fixed regardless of the number of FFN blocks.

3.2 Main Components

Base Model

The base model provides the trainable parameters within the MoFE architecture, including the embedding and self-attention layers of the entire architecture. TinyLlama, employed as the base model in the following experiments, features an embedding size of (32000, 2048) and 22 attention layers. In the MoFE architecture, the parameters provided by the base model are updated in contrast to the FFN blocks which remain frozen during the entire training process.

Expert Model

The FFN layers in the MoE architecture are

provided from the expert models. These FFN layers, which follow the attention layers in the Transformer architecture (Vaswani et al., 2023), primarily serve to maintain the isotropy of token embeddings (Sonkar and Baraniuk, 2023). As the FFN layers of TinyLlama comprise 0.76 billion parameters, integrating one expert model adds 0.76 billion, rather than the entire 1.1 billion parameters. As the FFN layers are frozen in MoFE, only the parameters located before the FFN blocks, which include the embeddings and self-attention layers provided by the base model, and the router, are updated.

Router

The router, or gate, includes a linear layer that determines which FFN block to activate for each token at every time step. This research uses a common gating method that leverages hidden state representations of positive and negative prompts, assigned during model merging. Routing assigns scores to each expert via a single matrix multiplication, computing dot products between a vector and the model’s hidden states to select the top two experts. Positive prompts are averaged, and negative prompts are subtracted to identify vectors that maximize these dot products.

4 Empirical Analysis

4.1 Experimental Setting

The experiments are implemented using three NVIDIA A100 80GB GPUs. The hyperparameters are set as follows: batch size of 4, learning rate of $3e-5$ with a linear learning rate scheduler, gradient accumulation of 512, and weight decay of 0.01.

Model	Fine-tuning	Trainable Parameters	Training Time(hr)	MMLU	MedMCQA
Small	X			0.2441	0.2678
	Full	1.86B	14	0.3331	0.3554
	MoFE	0.34B	6	0.3163	0.3431
Medium	X			0.2443	0.2661
	Full	3.38B	19	0.3231	0.3648
	MoFE	0.34B	6	0.3255	0.3297
Large	X			0.2448	0.2680
	Full	6.42B	26	0.3243	0.3459
	MoFE	0.34B	6	0.3130	0.3514

Table 1: Performance on MMLU and MedMCQA when the FFN blocks are updated and frozen, compared to before fine-tuning. All frozen models, regardless of size, have only 0.34 billion trainable parameters.

4.2 What is the trade-off between efficiency and performance?

To assess the impact of freezing FFN blocks on performance, we build MoFE models in three different sizes using the Mixtral architecture outlined in Section 3, with TinyLlama serving as both the base and expert models. We construct three models: a small model with 2 experts, a medium model with 4 experts, and a large model with 8 experts. Each model size is instruction-tuned using datasets from two distinct domains: MMLU (Hendrycks et al., 2021) for the general domain, and MedMCQA (Pal et al., 2022) for the medical domain. Since the MedMCQA training dataset contains approximately 18K rows, we randomly sample 18K rows from the MMLU dataset to ensure a balanced representation of both domains. We then train the models and compare their performance when the FFN blocks are either frozen or updated. The task performances are evaluated using lm-evaluation-harness (Gao et al., 2024).

Table 1 shows the number of trainable parameters, training time, and performance on MMLU and MedMCQA for models of each size when fully fine-tuned versus fine-tuned with FFN blocks frozen, referred to as MoFE. When fully fine-tuning, the number of trainable parameters increases with the number of expert models. However, in MoFE, the number of trainable parameters remains constant regardless of the number of expert models. This results in a fixed training time for MoFE models, while training time increases with model size for models with fully updated FFN blocks. Notably, even for the small model with 2 expert models, MoFE requires less than half the training time compared to fully updating the model.

To better understand the impact of each fine-tuning method, we also evaluate model perfor-

mance before fine-tuning. Both approaches improve performance, with full fine-tuning generally outperforming MoFE. However, exceptions exist: MoFE surpasses full fine-tuning on MMLU for the medium model and on MedMCQA for the large model. These findings suggest that while MoFE is slightly less effective overall, it remains competitive, offering significant efficiency gains in trainable parameters and training time. Appendix A further shows performance does not consistently correlate with the number of updated FFN blocks.

4.3 How good is MoFE compared to other PEFT methods?

Although MoFE demonstrates greater efficiency than full fine-tuning, it is important to compare MoFE with other PEFT methods to validate its effectiveness as an alternative training approach for low-resource environments. To this end, we utilize the same three model sizes—small, medium, and large—to compare the resource requirements and performance of various PEFT methods, including LoRA, QLoRA, and DoRA.

Table 2 demonstrates that among the four fine-tuning methods, MoFE consistently achieves the best performance on both MMLU and MedMCQA across all three model sizes. Despite having the highest number of trainable parameters, MoFE requires the least training time. These findings indicate that freezing the FFN blocks of MoE models can be an efficient fine-tuning approach, outperforming other PEFT methods by minimizing training time while maintaining strong performance on downstream tasks. Training time is a critical consideration in real-world scenarios, as it directly impacts computational costs, which scales linearly with GPU usage time.

Model	Fine-tuning	Trainable Parameters	Training Time(hr)	MMLU	MedMCQA
Small	X			0.2441	0.2678
	LoRA	2.3M	13	0.2935	0.2838
	QLoRA	2.3M	14	0.2953	0.2525
	DoRA	2.4M	15	0.2970	0.2682
	MoFE	0.34B	6	0.3163	0.3431
Medium	X			0.2443	0.2661
	LoRA	2.3M	15	0.2836	0.3053
	QLoRA	2.3M	15	0.2972	0.2608
	DoRA	2.4M	17	0.2934	0.3148
	MoFE	0.34B	6	0.3255	0.3297
Large	X			0.2448	0.2680
	LoRA	2.3M	18	0.2754	0.3091
	QLoRA	2.3M	22	0.2909	0.2682
	DoRA	2.4M	21	0.2935	0.2639
	MoFE	0.34B	6	0.3130	0.3514

Table 2: The number of trainable parameters, training time required, and performance on MMLU and MedMCQA using various fine-tuning methods. MoFE requires the least training time and achieves the best performance.

4.4 What effect does the domain expertise of consisting models have?

The MoFE architecture consists of two types of models: a base model and expert models, raising a key research question: How does the domain expertise of these models influence the overall performance of the MoFE model? To investigate this, we conduct a series of experiments focused on knowledge transfer from the consisting models.

4.4.1 Expert Model

Single Domain

To assess the impact of domain-specific knowledge in expert models, we build two separate models using TinyLlama: one trained on the MedMCQA dataset (medical expert model) and the other on the MMLU dataset (general model). We then construct several medium-sized MoFE models, each incorporating four expert models, where each expert is either a medical expert model or a general model. By varying the composition of these expert models, we aim to examine whether domain-specific knowledge from the expert models transfers to the overall MoFE model, with a particular focus on the medical domain. Since this experiment focuses on the impact of medical expert models, the base model is kept fixed as a general model without domain-specific expertise.

As shown in Table 3, performance on MedMCQA improves as the number of medical expert

Model		MedMCQA
Medical Expert	General	
0	4	0.3488
2	2	0.3536
4	0	0.3636

Table 3: The performance of MoFE models with various expert model compositions.

models increases. The model with four medical expert models achieves the highest performance, while the model without any medical expert models performs the lowest. This suggests that the presence of domain-specific expert models positively impacts the overall performance of the MoFE model, indicating that knowledge transfer from the expert models—specifically the FFN blocks—occurs within the MoFE architecture.

Multi-Domain

Building on the previous experiment confirming knowledge transfer in the medical domain, we investigate whether knowledge transfer across multiple domains is possible and how the number of domain-specific expert models affects the MoFE model’s domain knowledge. For this, we develop a finance expert model by training TinyLlama on the Sujet-Finance-Instruct-177k dataset¹, split

¹<https://sujet.ai/>

Model			Task Performance	
Finance Expert	Medical Expert	General	Medicine	Finance
0	0	4	0.3488	0.9087
0	2	2	0.3536	0.9237
0	4	0	0.3636	0.928
3	1	0	0.3603	0.936
2	2	0	0.3764	0.9327
1	3	0	0.3717	0.9401

Table 4: The performance of MoFE models with different numbers of finance expert models and medical expert models incorporated.

Base Model	Task Performance	
	Medicine	Finance
General	0.3763	0.9327
Medical Expert	0.3698	0.9326
Finance Expert	0.3598	0.9417

Table 5: Task performance of the MoFE models with different base models.

9:1 for training and testing. We then construct medium-sized MoFE models with varying numbers of medical expert models and finance expert models and evaluate them on MedMCQA and the Sujet-Finance-Instruct-177k test set. Finally, we compare these models with those from the **Single Domain** Section across both tasks.

As shown in Table 4, the MoFE model with two finance expert models and two medical expert models achieves the highest performance on MedMCQA, while the model with one finance expert model and three medical expert models performs best on Sujet-Finance-Instruct-177k. These findings suggest two key insights: incorporating domain-specific expert models enhances domain knowledge and task performance, but the number of domain expert models does not necessarily predict or linearly improve performance.

4.4.2 Base Model

The MoFE architecture requires not only expert models but also a base model that provides layers other than the FFN blocks, raising an additional research question: What is the impact of the base model’s domain expertise? Since our previous findings showed that including at least one domain expert model is crucial for domain-specific performance, we aim to isolate the influence of the base model in this experiment. To do so, we build three medium-sized MoFE models, each with a dif-

ferent base model: a general model, a medical expert model, and a finance expert model, while keeping the expert composition constant with two medical expert models and two finance expert models. We then evaluate these models on both medical and finance tasks.

As shown in Table 5, the MoFE model with the general model as the base performs best on the medical task and second best on the finance task. This suggests that using a general model as the base is a reasonable choice when building a MoFE model aimed at expertise across multiple domains.

4.5 What is the optimal training strategy?

Building on earlier experiments that demonstrated the potential for creating domain-specific expertise in MoFE models by incorporating pre-existing expert models, the next step is to determine the optimal training strategy for maximizing downstream task performance. Unlike prior experiments using only instruction-tuning, this section explores post-pretraining, where a pretrained model undergoes additional pretraining before fine-tuning. The goal is to assess whether a pretrained or instruction-tuned model as the expert is more effective and if post-pretraining adds value or instruction-tuning alone is sufficient for MoFE models.

For testing in the medical domain, the pretraining datasets include English data from the Multilingual-Medical-Corpus (García-Ferrero et al., 2024) for the medical domain and Multi-News data (Fabbri et al., 2019) for the general domain. Due to dataset distribution balance, a random sample of 0.2 million rows from each dataset, totaling 0.4 million rows, is used for training. The MedMCQA instruction dataset is used for instruction-tuning across all strategies. Medium-sized MoFE models with four expert models are tested under the following training strategies:

Expert Model	Training Strategy	MedMCQA	PubMedQA
TinyLlama	Instruction-tuning	0.3529	0.6
	Post-pretraining → Instruction-tuning	0.2589	0.188
Medical Expert	Instruction-tuning	0.3655	0.584

Table 6: Task performance across various training strategies.

1. Using TinyLlama, as the expert models, followed by instruction-tuning the MoFE model.
2. Using TinyLlama as the expert models, post-pretraining, and then instruction-tuning the MoFE model.
3. Using the medical expert model, as the expert models, followed by instruction-tuning the MoFE model.

For evaluation, we use two medical tasks: MedMCQA and PubMedQA (Jin et al., 2019). PubMedQA, derived from PubMed abstracts², serves as an additional benchmark since the medical expert models were trained with MedMCQA data, which could inflate performance by resembling additional training epochs. To ensure a fairer comparison, we evaluate the models on PubMedQA, an unseen dataset, to test medical knowledge.

We compare MoFE models using TinyLlama as expert models under both instruction-tuning alone and post-pretraining followed by instruction-tuning, but only test the MoFE model with medical expert models under instruction-tuning. This is because the medical expert models are already instruction-tuned, and post-pretraining an instruction-tuned model leads to catastrophic forgetting, reducing performance, as noted by Luo et al. (2024).

Table 6 shows that performance on both MedMCQA and PubMedQA is worst with the second strategy, involving post-pretraining followed by instruction-tuning with TinyLlama as expert models. The best strategies differ: for MedMCQA, the third strategy, using medical expert models followed by instruction-tuning, is optimal, while for PubMedQA, the first strategy, using TinyLlama as expert models and instruction-tuning without post-pretraining, yields the best performance.

The superior performance of the third strategy for MedMCQA is expected, as the medical expert model is TinyLlama instruction-tuned with

MedMCQA data resulting in the same effect of undergoing an additional training epoch. Since PubMedQA is a completely unseen task, it serves as a more objective performance indicator. The results suggest that the first strategy, using TinyLlama as expert models and instruction-tuning the MoFE model directly, is the optimal approach.

The results indicate that post-pretraining significantly decreases performance on both tasks, which can be explained by the characteristics of the MoFE architecture. Integrating new knowledge effectively requires updating all layers of the model, but FFN blocks remain frozen in MoFE. Given that FFN layers constitute a significant portion of the model’s parameters, they likely play a crucial role in knowledge integration. Language models primarily acquire knowledge during pretraining, with instruction-tuning focused on adapting to specific task formats rather than acquiring new knowledge (Zhao et al., 2023). Therefore, post-pretraining a model with frozen FFN layers, where only the parameters before these layers are updated, may result in misalignment among the various model layers. This misalignment could possibly explain the observed decrease in performance when using post-pretraining.

5 Conclusion

Given the enormous computational costs of training and serving LLMs, we propose MoFE as an efficient model training and scaling strategy. While there is a trade-off between efficiency and performance, MoFE significantly reduces the size of trainable parameters and training time, demonstrating superiority over other PEFT methods in both training time and task performance. Furthermore, the transfer of domain expertise from the constituent models enables the creation of multi-domain proficient models by leveraging existing domain experts. We believe MoFE presents a viable option for resource-constrained environments in real-world scenarios.

²<https://pubmed.ncbi.nlm.nih.gov/>

Limitation

The base and expert models used in this work is relatively small, with only 1.1 billion parameters. For lightweight experiments, we utilized a limited amount of data from a few domains. Consequently, the experimental results cannot be fully generalized to larger models or all domains.

Ethics Statement

Given that computational costs entail not only monetary issues but also environmental concerns, we strive to provide as much information as possible to facilitate the reproduction of our experiments. Further, although we refer to the models instruction-tuned with medical data and finance data as medical expert model and finance expert model respectively, these names are for simplicity in reference only. These models should not be considered actual domain experts capable of providing clinical or financial advice.

References

- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. [Qlora: Efficient finetuning of quantized llms](#). *Preprint*, arXiv:2305.14314.
- Alexander Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir Radev. 2019. [Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1074–1084, Florence, Italy. Association for Computational Linguistics.
- William Fedus, Barret Zoph, and Noam Shazeer. 2022. [Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity](#). *Preprint*, arXiv:2101.03961.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2024. [A framework for few-shot language model evaluation](#).
- Iker García-Ferrero, Rodrigo Agerri, Aitziber Atutxa Salazar, Elena Cabrio, Iker de la Iglesia, Alberto Lavelli, Bernardo Magnini, Benjamin Molinet, Johana Ramirez-Romero, German Rigau, Jose Maria Villa-Gonzalez, Serena Villata, and Andrea Zaninello. 2024. [Medical mt5: An open-source multilingual text-to-text llm for the medical domain](#). *Preprint*, arXiv:2404.07613.
- Charles Goddard, Shamane Siriwardhana, Malikeh Ehghaghi, Luke Meyers, Vlad Karpukhin, Brian Benedict, Mark McQuade, and Jacob Solawetz. 2024. [Arcee’s mergekit: A toolkit for merging large language models](#). *Preprint*, arXiv:2403.13257.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). *Preprint*, arXiv:2009.03300.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. [Parameter-efficient transfer learning for nlp](#). *Preprint*, arXiv:1902.00751.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#). *Preprint*, arXiv:2106.09685.
- Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. 2017. [Quantization and training of neural networks for efficient integer-arithmetic-only inference](#). *Preprint*, arXiv:1712.05877.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *Preprint*, arXiv:2310.06825.
- 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, L lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. 2024. [Mixtral of experts](#). *Preprint*, arXiv:2401.04088.
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. [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, Hong Kong, China. Association for Computational Linguistics.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. [Scaling laws for neural language models](#). *Preprint*, arXiv:2001.08361.

- Aran Komatsuzaki, Joan Puigcerver, James Lee-Thorp, Carlos Riquelme Ruiz, Basil Mustafa, Joshua Ainslie, Yi Tay, Mostafa Dehghani, and Neil Houlsby. 2023. [Sparse upcycling: Training mixture-of-experts from dense checkpoints](#). *Preprint*, arXiv:2212.05055.
- Yanis Labrak, Adrien Bazoge, Emmanuel Morin, Pierre-Antoine Gourraud, Mickael Rouvier, and Richard Dufour. 2024. [Biomistral: A collection of open-source pretrained large language models for medical domains](#). *Preprint*, arXiv:2402.10373.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). *Preprint*, arXiv:2104.08691.
- Opher Lieber, Barak Lenz, Hofit Bata, Gal Cohen, Jhonathan Osin, Itay Dalmedigos, Erez Safahi, Shaked Meirom, Yonatan Belinkov, Shai Shalev-Shwartz, Omri Abend, Raz Alon, Tomer Asida, Amir Bergman, Roman Glozman, Michael Gokhman, Avshalom Manevich, Nir Ratner, Noam Rozen, Erez Shwartz, Mor Zusman, and Yoav Shoham. 2024. [Jamba: A hybrid transformer-mamba language model](#). *Preprint*, arXiv:2403.19887.
- Shih-Yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting Cheng, and Min-Hung Chen. 2024. [Dora: Weight-decomposed low-rank adaptation](#). *Preprint*, arXiv:2402.09353.
- Yun Luo, Zhen Yang, Fandong Meng, Yafu Li, Jie Zhou, and Yue Zhang. 2024. [An empirical study of catastrophic forgetting in large language models during continual fine-tuning](#). *Preprint*, arXiv:2308.08747.
- Ankit Pal, Logesh Kumar Umapathi, and Malaikannan Sankarasubbu. 2022. [Medmcqa : A large-scale multi-subject multi-choice dataset for medical domain question answering](#). *Preprint*, arXiv:2203.14371.
- Sara Pieri, Sahal Shaji Mullappilly, Fahad Shahbaz Khan, Rao Muhammad Anwer, Salman Khan, Timothy Baldwin, and Hisham Cholakkal. 2024. [Bimedix: Bilingual medical mixture of experts llm](#). *Preprint*, arXiv:2402.13253.
- Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarczyk, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. [Outrageously large neural networks: The sparsely-gated mixture-of-experts layer](#). *Preprint*, arXiv:1701.06538.
- Shashank Sonkar and Richard G. Baraniuk. 2023. [Investigating the role of feed-forward networks in transformers using parallel attention and feed-forward net design](#). *Preprint*, arXiv:2305.13297.
- Katrin Tomanek, Vicky Zayats, Dirk Padfield, Kara Vailancourt, and Fadi Biadisy. 2021. [Residual adapters for parameter-efficient ASR adaptation to atypical and accented speech](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6751–6760, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open foundation and fine-tuned chat models](#). *Preprint*, arXiv:2307.09288.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. [Attention is all you need](#). *Preprint*, arXiv:1706.03762.
- Fuzhao Xue, Zian Zheng, Yao Fu, Jinjie Ni, Zangwei Zheng, Wangchunshu Zhou, and Yang You. 2024. [Openmoe: An early effort on open mixture-of-experts language models](#). *Preprint*, arXiv:2402.01739.
- Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024. [Tinyllama: An open-source small language model](#). *Preprint*, arXiv:2401.02385.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. [A survey of large language models](#). *Preprint*, arXiv:2303.18223.

A Does the number of frozen FFN blocks affect performance?

FFN Blocks		MedMCQA
Frozen	Updated	
4	0	0.3529
3	1	0.3407
2	2	0.3524
1	3	0.3541
0	4	0.3705

Table 7: The effect of the number of frozen FFN blocks on task performance.

To examine how performance shifts with varying numbers of frozen FFN blocks, we use a medium-sized model with four expert models. TinyLlama serves as the base and expert models, as in previous experiments. Five versions of the model are constructed: one with all expert models frozen, one with three frozen, one with two frozen, one with one frozen, and one with none frozen. Each model is instruction-tuned on the MedMCQA training dataset and evaluated on its test set.

As shown in Table 7, the fully updated model demonstrated the best performance. However, the results reveal that performance does not consistently correlate with the number of frozen FFN blocks as expected.