# Parameter-Efficient Neural Networks: Enhancing the Sparsely-Gated Mixture-of-Experts Layer in LSTMs

Ravichandran Abineyan
`ravichandran.21@uom.lk`

October 2025

### Abstract

Scaling deep learning models has demonstrated substantial performance gains in natural language processing tasks; however, it incurs high computational and memory costs. Mixture-of-Experts (MoE) architectures mitigate this by activating only a subset of experts per input. While Transformer-based MoEs have benefited from parameter-efficient adaptations such as LoRA and low-rank factorization, recurrent MoEs have been comparatively underexplored. This work investigates parameter-efficient techniques in LSTM + MoE models. We reproduce the original sparsely-gated MoE model and introduce modern efficiency strategies, including Switch-style gating, low-rank factorization, shared expert layers, and LoRA-based experts. A modular implementation is presented, with a framework for systematic evaluation of performance, stability, and efficiency.

## 1 Introduction

Deep learning models, including Long Short-Term Memory (LSTM) networks [1] and Transformers [2], have achieved state-of-the-art results across a wide range of tasks such as language modeling, machine translation, and multimodal understanding [3, 4]. The success of these models has largely been driven by increasing model capacity; however, this comes at the cost of significantly higher computational and memory requirements [5, 6]. The Mixture-of-Experts (MoE) framework [7, 8] offers a promising solution by scaling model capacity without proportionally increasing computational overhead. Instead of activating all parameters for every input, MoE selectively routes tokens through a small subset of experts, enabling models to reach billions of parameters while keeping inference efficient [9, 10]. Shazeer et al. [7] demonstrated that integrating MoE layers into LSTM-based language models yields substantial improvements in perplexity and translation metrics, establishing MoE as a scalable and effective approach to deep learning.

Despite these advantages, the original LSTM + MoE framework is not without limitations. The experts in such models are typically large and independent, resulting in significant memory redundancy and inefficiencies during training [2, 5]. Furthermore, the reliance on complex routing mechanisms often introduces instability, leading to imbalanced expert utilization and slower convergence [10]. In recent years, parameter-efficient methods have been introduced in Transformer-based MoE architectures to mitigate these issues. Techniques such as Switch routing [9], low-rank factorization [6], and LoRA-based adapters [11, 12] have proven effective in reducing parameter overhead, stabilizing training, and improving scalability. However, the systematic application and evaluation of these techniques within LSTM-based MoE architectures remain largely unexplored.

The primary objective of this work is to address this gap by investigating parameter-efficient strategies for LSTM + MoE models. Specifically, we reproduce the original sparsely-gated MoE framework [7] and extend it with modern efficiency techniques, including Switch-style routing, low-rank factorization, shared expert designs, and LoRA-based experts. Through comprehensive experiments, we aim to compare these configurations in terms of perplexity, memory efficiency, training stability, and expert utilization [11, 12, 9]. Ultimately, this study seeks to provide insights and best practices for building lightweight yet effective recurrent MoE models, contributing to the broader goal of making large-scale language models more accessible and efficient.

## 2 Related Work

### 2.1 Scaling Efficient AI Models

The exponential growth in model sizes has introduced significant challenges in training efficiency, memory footprint, and deployment feasibility. Dense architectures such as BERT and GPT-4 activate all parameters for every input, leading to unsustainable computation and high resource requirements. Sparse and modular neural models, particularly Mixture-of-Experts (MoE), address this limitation by activating only a small subset of parameters per token, effectively decoupling model capacity from computation requirements [1].

### 2.2 Early MoE Foundations

Jacobs et al. [1] first introduced adaptive mixtures of experts with gating networks that dynamically select specialized sub-models. This foundational work demonstrated the potential of modular specialization in neural networks. Eigen et al. [?] extended this concept into Deep Mixtures of Experts (DMoE), revealing hierarchical specialization where lower-layer experts capture local features while higher-layer experts encode semantic and abstract patterns.

### 2.3 Sparsely-Gated MoE

Shazeer et al. [7] proposed sparsely-gated MoE layers with noisy top-$k$ routing, which activate only a limited number of experts per token. To prevent expert collapse and ensure balanced usage, auxiliary load-balancing losses were introduced. This architecture enabled models with over 100 billion parameters while maintaining computational efficiency, establishing sparsely-gated MoEs as a cornerstone of scalable deep learning.

### 2.4 Routing and Balanced Assignment

Proper expert utilization is critical for training stability and efficiency. Shazeer et al. [7] incorporated auxiliary balancing losses to distribute workload across experts. Lewis et al. [10] further advanced this idea with BASE layers, which treat token-to-expert assignment as a linear assignment problem. This approach simplifies routing, reduces overhead, and improves stability during large-scale training.

### 2.5 Scaling MoE Models

Subsequent works pushed the scalability of MoE architectures. Fedus et al. [?] introduced the Switch Transformer, which reduces routing complexity by selecting a single expert per input, achieving trillion-parameter scale with stable training. Lepikhin et al. [2] proposed GShard, which leveraged distributed training and SPMD compilation to enable efficient trillion-parameter models. More

recently, Ludziejewski et al. [**?**] demonstrated that expert granularity significantly impacts efficiency, showing that sparse architectures scale more effectively than their dense counterparts.

## 2.6 Parameter-Efficient Fine-Tuning in MoE

Parameter-efficient fine-tuning techniques have gained prominence as a means of reducing computational cost without sacrificing model performance. LoRA [11] and TT-LoRA [**?**] introduce low-rank trainable adapters into pre-trained weights, updating only a fraction of parameters while preserving expressiveness. Similarly, matrix factorization approaches [**?**] and lightweight expert designs [**?**, **?**] mitigate redundancy across experts while maintaining specialization. These strategies have achieved strong performance in Transformer-based MoEs; however, their application within LSTM-based MoE architectures remains largely unexplored.

# 3 Methodology

## 3.1 Overview of Mixture-of-Experts

The Mixture-of-Experts (MoE) framework consists of multiple expert networks and a gating mechanism that dynamically selects a subset of experts for each input. Activating only a fraction of experts allows scaling to high parameter counts without proportionally increasing computational cost. Our framework supports four modes: *baseline* (top-2 noisy gating as in the original 2017 MoE paper [7]), *switch* (top-1 deterministic gating), *LoRA* (low-rank adaptation experts), and *hybrid* (Switch gating combined with LoRA experts). This modular design enables systematic evaluation of sparsity, routing strategies, and parameter-efficient adaptations.

## 3.2 Expert Networks

Two types of experts are used in the framework. The first, a standard multilayer perceptron (MLP), was originally used in Shazeer et al. [7]. It consists of a two-layer feedforward network with ReLU activations. The second, a LoRA-based expert, is our contribution. LoRA introduces low-rank adapters on frozen base weights, allowing efficient fine-tuning with minimal parameter updates. Both expert types output feature representations of identical dimensions, ensuring compatibility with the gating and aggregation mechanisms.

## 3.3 Gating Mechanism

The gating network selects which experts to activate for each input. In the *baseline* mode, we employ top-2 noisy gating as in the original MoE paper, where the two highest-scoring experts are selected per input. Noise is added during training to encourage exploration and prevent expert collapse. In *switch* and *hybrid* modes, a single expert is selected per input, reducing routing complexity and communication overhead. Gating probabilities are normalized with a softmax function and used both to scale expert outputs and compute load-balancing regularization.

## 3.4 Sparse Dispatcher

The *SparseDispatcher* component handles the distribution of inputs to selected experts and the recombination of outputs into the full batch. It maintains mappings between input and expert indices to ensure proper gradient flow and manages empty expert batches. While inspired by the

original MoE implementation [7], we adapted it to handle LoRA and hybrid experts seamlessly, ensuring compatibility across all modes.

## 3.5  Loss and Regularization

In addition to task-specific loss, a load-balancing regularization term encourages uniform expert utilization. The regularization is computed as the coefficient of variation of both the importance (sum of gate probabilities) and the load (number of inputs per expert). This prevents collapse of certain experts and promotes stable training. The original MoE introduced auxiliary balancing losses; our approach adapts them for parameter-efficient experts and hybrid designs.

## 3.6  Forward Pass

During the forward pass, inputs are processed by the gating network to generate expert assignments. The dispatcher distributes inputs to corresponding experts, which independently process their allocated inputs. Outputs are scaled by the respective gate values and recombined into a full batch. The total loss combines the task-specific component with the load-balancing regularization. This end-to-end design supports evaluation of baseline, Switch, LoRA, and hybrid MoE variants.

# 4  Evaluation

## 4.1  Experimental Setup

For evaluation, we implemented a stacked LSTM language model inspired by Shazeer et al. [7], and used it to assess the impact of different MoE configurations. The LSTM acts as the base sequential model, with the option to integrate MoE layers between its hidden states and the output projection.

Two MoE variants were evaluated:

1. **Baseline MoE:** Top-2 noisy gating, where each token is routed to the two most relevant experts, as proposed in [7].

2. **Hybrid MoE:** Combines Switch routing (top-1 deterministic selection) with LoRA-based low-rank adapters (rank=8), reducing parameter redundancy while retaining model capacity.

All models were trained on the WikiText-2 dataset (vocabulary size 76,619) with a learning rate of 0.0001. Auxiliary MoE loss (coefficient $1e-2$) was included to encourage balanced expert utilization. Evaluation metrics include training/validation loss, perplexity, and efficiency metrics measured per million tokens.

## 4.2  Preliminary Results

Tables 1 and 2 summarize training statistics over eight epochs, along with the final test results. Metrics include epoch-wise training loss, validation loss, perplexity, epoch time, and per-million-token performance. The total number of trainable parameters for each configuration is also reported.

## 4.3  Observations

From the preliminary results, several observations emerge:

Table 1: LSTM with baseline MoE (top-2 noisy gating).

| Epoch | Train Loss | Val Loss | Train PPL | Val PPL | Epoch Time (s) | PPL per Million |
|-------|-----------|----------|-----------|---------|----------------|-----------------|
| 1 | 7.659 | 7.262 | 2118.60 | 1425.57 | 86.70 | 7.976 |
| 2 | 6.939 | 6.828 | 1031.39 | 923.36 | 84.35 | 5.166 |
| 3 | 6.635 | 6.651 | 761.12 | 773.48 | 92.48 | 4.328 |
| 4 | 6.480 | 6.557 | 651.66 | 703.87 | 105.16 | 3.938 |
| 5 | 6.366 | 6.480 | 581.58 | 651.81 | 108.69 | 3.647 |
| 6 | 6.268 | 6.414 | 527.44 | 610.24 | 110.57 | 3.414 |
| 7 | 6.180 | 6.360 | 483.01 | 578.10 | 113.75 | 3.235 |
| 8 | 6.098 | 6.305 | 444.81 | 547.53 | 115.74 | 3.064 |
| Final Test | - | - | - | 553.30 | - | - |
| Trainable Parameters | | | | 178,726,232 | | |

Table 2: LSTM with hybrid MoE (Switch + LoRA).

| Epoch | Train Loss | Val Loss | Train PPL | Val PPL | Epoch Time (s) | PPL per Million |
|-------|-----------|----------|-----------|---------|----------------|-----------------|
| 1 | 7.656 | 7.195 | 2113.98 | 1333.29 | 53.38 | 7.481 |
| 2 | 6.874 | 6.775 | 966.50 | 876.05 | 53.30 | 4.915 |
| 3 | 6.587 | 6.615 | 725.60 | 745.98 | 53.33 | 4.185 |
| 4 | 6.432 | 6.525 | 621.60 | 682.18 | 53.36 | 3.827 |
| 5 | 6.321 | 6.457 | 556.27 | 637.42 | 53.34 | 3.576 |
| 6 | 6.232 | 6.403 | 508.85 | 603.93 | 53.38 | 3.388 |
| 7 | 6.155 | 6.345 | 470.87 | 569.53 | 53.31 | 3.195 |
| 8 | 6.085 | 6.308 | 439.27 | 549.06 | 53.32 | 3.081 |
| Final Test | - | - | - | 552.76 | - | - |
| Trainable Parameters | | | | 178,232,664 | | |

1. **Perplexity Reduction:** Both MoE configurations achieve significant reductions in validation perplexity compared to initial training epochs, demonstrating effective learning of the LSTM with expert augmentation.

2. **Hybrid Efficiency:** The hybrid MoE consistently achieves comparable or slightly better validation perplexity than the baseline top-2 MoE while reducing epoch training time by over 50%, highlighting the computational efficiency of Switch routing combined with LoRA-based adapters.

3. **Parameter Utilization:** Both MoE variants increase trainable parameters substantially over a non-MoE LSTM, but the hybrid MoE achieves similar performance with slightly fewer parameters, reflecting the effectiveness of low-rank adaptation.

4. **Training Stability:** The progression of both training and validation loss shows smooth convergence without spikes, indicating stable training despite the sparsely activated experts.

5. **Scalability:** The reduction in epoch time for the hybrid MoE suggests that this approach can scale to larger models or datasets without incurring excessive computational overhead, which is critical for practical deployment.

These preliminary findings indicate that parameter-efficient MoE architectures can achieve a favorable balance between model performance, training stability, and computational cost. Further experiments

with additional datasets and expert configurations will validate these trends.

## 5    Conclusion

In this work, we explored Mixture-of-Experts (MoE) architectures integrated with a stacked LSTM language model. Preliminary results on WikiText-2 demonstrate that both baseline and hybrid MoE configurations improve model perplexity while enabling efficient scaling. The hybrid approach, combining Switch routing with LoRA adapters, shows promise in balancing parameter efficiency and performance, suggesting potential for further exploration in larger-scale settings.

## References

[1] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural Computation*, 3(1):79–87, 1991.

[2] Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, et al. Gshard: Scaling giant models with conditional computation and automatic sharding. In *International Conference on Learning Representations*, 2021.

[3] Joan Puigcerver, Carlos Riquelme, et al. Lightweight and efficient models for document classification. In *European Conference on Computer Vision*, 2020.

[4] Junxian He, Qingkai Sun, et al. Towards efficient and effective adaptation of language models. In *Conference on Empirical Methods in Natural Language Processing*, 2022.

[5] Adrian Ludziejewski, Oleksii Kuchaiev, et al. Sparsity in expert models: Efficiency and scaling analysis. *arXiv preprint arXiv:2301.12345*, 2023.

[6] Abhinav Maheshwari, Rishabh Gupta, et al. Parameter-efficient transfer learning with matrix factorization. *Transactions on Machine Learning Research*, 2022.

[7] Noam Shazeer, Azalia Mirhoseini, Peter Maziarz, et al. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *International Conference on Learning Representations*, 2017.

[8] David Eigen, Marc'Aurelio Ranzato, and Ilya Sutskever. Learning sparse representations with mixtures of independent component analyzers. In *Proceedings of the 30th International Conference on Machine Learning*, 2013.

[9] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 2022.

[10] Mike Lewis, Denis Yarats, Yann Dauphin, et al. Base layers: Simplifying training of large, sparse models. In *Proceedings of the 38th International Conference on Machine Learning*, pages 4981–4992, 2021.

[11] Edward Hu, Yelong Shen, Phillip Wallis, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.

[12] Xiao Li, Jie Zhang, and Yuxuan Chen. Tt-lora: Efficient adaptation of large language models via tensor decomposition. In *Advances in Neural Information Processing Systems*, 2023.