

# MIXTURE OF NEURON EXPERTS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In this work, We first explore whether the parameters activated by the MoE layer remain highly sparse at inference. We perform a sparsification study on several representative MoE models. For each expert, we rank parameters by the magnitude of their activations from the gate projection and progressively prune the activated subset. Pruning up to 60% of parameters within that subset causes only negligible task-performance degradation; substantial drops occur only after more than 90% are removed. We further decompose experts into neuron granular MoE and visualize their activation values, finding that most neuron activations are near zero. This observation motivates us to select only high-activation neuron experts during pretraining. Based on this insight, we propose *Mixture of Neuron Experts* (MoNE). MoNE achieve neuron granular expert selection by only applying a simple top- $k$  selection within each expert, incurs negligible latency, and requires no additional routing parameters or inter-expert communication. Extensive experiments demonstrate that MoNE matches traditional MoE performance while activating only 50% of the MoE-layer parameters, and it consistently outperforms traditional MoE when compared at equal numbers of activated parameters. These results suggest that MoNE is a practical approach to improving parameter utilization and inference efficiency in MoE-like models.

## 1 INTRODUCTION

Large language models (LLMs) (Dai et al., 2024; Bai et al., 2023; Agarwal et al., 2025; Team et al., 2025) based on Mixture-of-Experts approaches have attracted growing interest in both academic research and industry. The fundamental concept of Mixture-of-Experts (MoE) in large language models entails partitioning a large feed-forward network (FFN) into several smaller subnetworks referred to as experts, where only a subset of expert parameters are activated depending on the input. Unlike dense models that activate all parameters uniformly, MoE models achieve greater computational efficiency through sparse activation patterns.

One key motivation for MoE is the long-observed activation sparsity (Frankle & Carbin, 2018; Fedus et al., 2022a; Frantar & Alistarh, 2023; Frankle et al., 2019) in dense networks: for an given input, only a small fraction of parameters are effectively used. Mixture-of-Experts architectures exploit this property via conditional computation, activating a sparse subset of specialists so as to substantially increase model capacity while maintaining computational efficiency. This further motivates us to propose the following question:

*Are the parameters activated by the MoE layer still highly sparse at inference?*

To answer the question, we performed a sparsification study on a set of representative MoE models. For each expert, we ranked the parameters according to the magnitude of their activations values, which calculated by the gate projection. Then we progressively pruned the weights from the activated subset according to their rank. The results are presented in Figure 1. Notably, across three evaluated models, removing up to 60% of the parameters in this subset led to only negligible declines in task performance, with significant degradation occurring only after more than 90% were pruned. These results suggest that the parameter subset selected by the MoE gating mechanism still highly sparsity at inference.

To further explore the sparsity of MoE, we decompose expert into neuron granular MoE. Then we visualize the activation value for the neuron experts. As shown in Figure 2, most of the activation

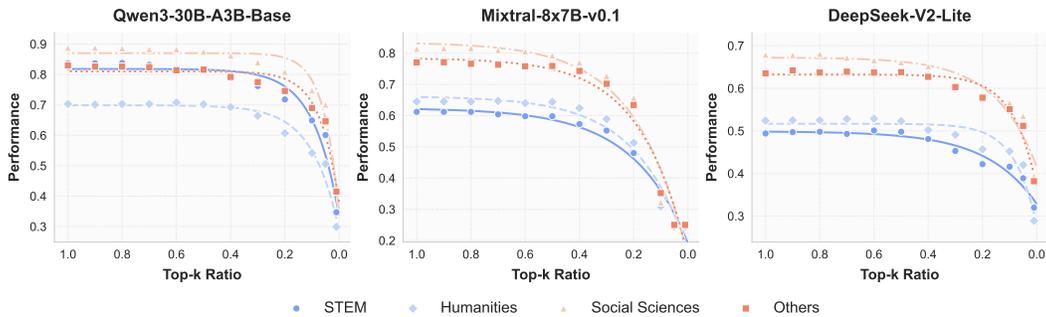


Figure 1: The performance of mainstream MoE models when only use the neuron experts with higher activation weight without extra training. Top-K Ratio refers to the ratio of selected neuron experts.

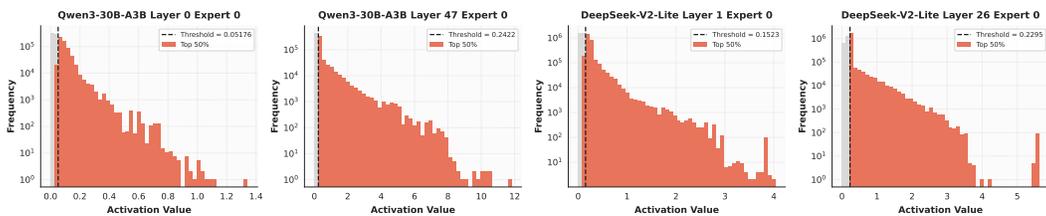


Figure 2: The activation value for the neuron experts, and the top 50% of these values were highlighted.

values are small, which further demonstrate the sparsity of the MoE layer. This results motivate us to only use the neuron experts with high activation weights for pretraining. Also, recently studies have shown the importance of expert granularity (Krajewski et al., 2024; Lepikhin et al., 2020; Du et al., 2022): Deepseek V3 Liu et al. (2024) applies 256 experts, Kimi K2 Team et al. (2025) applies 384 experts, and Qwen3-NextTeam (2025) applies 512 experts. However, overly fine-grained expert partitioning requires substantially larger routing networks and incurs significant cross-device communication latency (Lepikhin et al., 2020; Fedus et al., 2022b). Therefore, we propose Mixture of Neuron Experts (MoNE). By applying a simple top-k selection within each expert, we achieve granular selection for MoE without introducing additional router parameters or inter-expert communication. We evaluate MoNE under multiple settings and find that it consistently outperforms traditional MoE.

Our contributions are summarized as follows:

- We empirically show that traditionally trained Mixture-of-Experts models exhibit high activation sparsity at inference: a small subset of parameters with large activation values retains most of the model’s capability, and most of the activation values are small in the MoE layer.
- We introduce *Mixture of Neuron Experts* (MoNE), which decompose the expert into neuron granular MoE, and achieve neuron granular expert selection via a simple top-k within-expert operation that incurs negligible latency overhead and requires no additional routing parameters.
- Extensive experiments demonstrate that MoNE matches the performance of traditional MoE while using only 50% of the parameters in MoE layer. With the same total number of activated parameters, MoNE consistently outperforms traditional MoE.

## 2 RELATED WORK

### 2.1 LARGE LANGUAGE MODELS

Large language models (LLM) (Touvron et al., 2023a; Bai et al., 2023; Brown et al., 2020; Achiam et al., 2023; Liu et al., 2024; Devlin et al., 2019; Raffel et al., 2020) have shown remarkable abil-

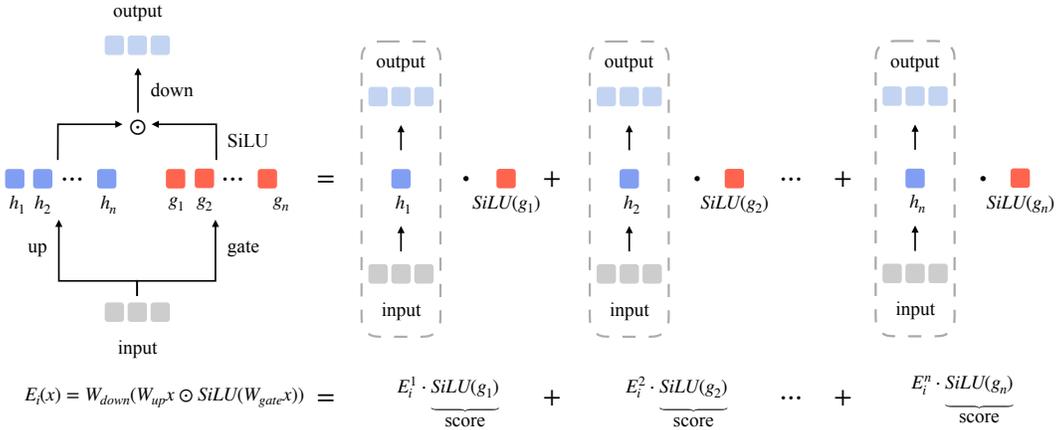


Figure 3: Expert in traditional MoE can be decomposed as the weighted sum of neuron granular FFN, which can be realized as a neuron granular MoE.

ities across different tasks, representing important progress toward artificial general intelligence. This success is largely driven by the growth of training data and the expansion of model parameter counts (Wei et al., 2022; Kaplan et al., 2020). And many recent works successfully scaling LLM to billions of parameters (Dai et al., 2024; Liu et al., 2024; Team et al., 2025; Agarwal et al., 2025; Zhang et al., 2022; Scao et al., 2022). However, As model scale increases, the demand for computational resources rises sharply. Consequently, improving the efficiency of both training and inference has become a central research focus to enable further scaling of large language models.

## 2.2 MIXTURE OF EXPERTS

The Mixture of Experts (MoE) (Cai et al., 2025; Masoudnia & Ebrahimpour, 2014; Jiang et al., 2024) architecture was introduced to enhance the capacity of deep neural networks while maintaining computational efficiency. Shazeer et al. (2017) proposed integrating an MoE layer between LSTM layers, demonstrating strong performance in language modeling and machine translation tasks. This approach was later adapted into the transformer framework by replacing the standard feed-forward layers with MoE modules. The Switch Transformer (Fedus et al., 2022b) streamlines the expert selection process by assigning each token to only the top-ranked expert, enabling more efficient model scaling. Gshard (Lepikhin et al., 2020) refined the routing mechanism by employing a Top-2 expert strategy, leading to substantial improvements in multilingual translation across 100 languages. More recently, DeepseekMoE (Dai et al., 2024) and have introduced fine-grained partitioning of experts within the MoE structure. Grove-MoE (Wu et al., 2025) incorporating experts of varying sizes. PEER (He, 2024) scales the number of experts up to one million. Kimi-K2 (Team et al., 2025) scales the model to 1,000B parameters and employs 384 experts, while Qwen3-Next (Team, 2025) further increases the expert count to 512. Improving expert granularity (Tian et al., 2025; Krajewski et al., 2024) and increasing the utilization of activated parameters (Li et al., 2025) are central goals in contemporary MoE research. However, overly fine-grained expert partitioning requires substantially larger routing networks and incurs significant cross-device communication overhead and latency (Lepikhin et al., 2020; Fedus et al., 2022b). PKM (Lample et al., 2019) achieves fast retrieval in a key-value memory layer using a structured memory architecture and further scaled by (Huang et al., 2024). In this work, we analyze the sparsity of activated parameters in traditional MoE architectures and construct a neuron granularity MoE model. This neuron granularity conversion substantially improves the utilization of activated parameters while avoiding the need for a large router and communication latency associated with overly fine expert partitioning in traditional MoE.

### 3 METHOD

#### 3.1 PRELIMINARY ON MOE

The Mixture-of-Experts (MoE) layer extends a standard Transformer by replacing a dense feed-forward network (FFN) with a collection of expert FFNs and a routing mechanism. For each input token, the router conditionally dispatches only a sparse subset of experts, so that computation is performed by a small number of specialists rather than the entire FFN. This sparse execution increases the model’s effective capacity while keeping per-token computation and latency approximately constant, enabling parameter-efficient scaling to much larger models. (Lepikhin et al., 2020; Dai et al., 2024; Fedus et al., 2022b).

Formally, let  $\mathbf{x} \in \mathbb{R}^{d_{\text{model}}}$  denote an input hidden state. An MoE layer consists of  $N_E$  experts  $\{E_i\}_{i=1}^{N_E}$  together with a router that maps  $\mathbf{x}$  to routing logits. The experts typically use the *Gated Linear Unit* (Dauphin et al., 2017) structure, which can be formulated as:

$$\mathbf{E}_i(\mathbf{x}) = \mathbf{W}_{\text{down}}^i (\text{SiLU}(\mathbf{W}_{\text{gate}}^i \mathbf{x}) \odot \mathbf{W}_{\text{up}}^i \mathbf{x}) \quad (1)$$

where  $\mathbf{W}_{\text{gate}}^i \in \mathbb{R}^{d_{\text{model}} \times d_{\text{expert}}}$  is the gate projection,  $\mathbf{W}_{\text{up}}^i \in \mathbb{R}^{d_{\text{model}} \times d_{\text{expert}}}$  is the up projection, and  $\mathbf{W}_{\text{down}}^i \in \mathbb{R}^{d_{\text{expert}} \times d_{\text{model}}}$  is the down projection. A common routing pipeline is to first calculate the scores of the router:

$$\mathbf{P}(\mathbf{x}) = \text{Act}(\text{topK}(\text{Router}(\mathbf{x}))), \quad (2)$$

where  $\text{topK}(\cdot)$  masks out all but the top-K routing logits and  $\text{Act}(\cdot)$  is the activation function. Then the MoE output is calculated as follows:

$$\text{MoE}(\mathbf{x}) = \sum_{i=1}^{N_E} \mathbf{P}(\mathbf{x})_i \mathbf{E}_i(\mathbf{x}). \quad (3)$$

To encourage balanced utilization across experts, we adopt the commonly used auxiliary load-balance loss:

$$\mathcal{L}_{\text{aux}} = \alpha_{\text{aux}} \cdot N_E \cdot \sum_{i=1}^{N_E} \mathbf{f}_i \cdot \mathbf{P}_i, \quad \text{where} \quad (4)$$

$$\mathbf{f}_i = \frac{1}{T} \sum_{\mathbf{x} \in \mathcal{B}} \mathbb{1}\{i \in \text{argtopK}(\text{Router}(\mathbf{x}))\}, \quad \mathbf{P}_i = \frac{1}{T} \sum_{\mathbf{x} \in \mathcal{B}} \frac{\text{Act}(\text{Router}(\mathbf{x}))[i]}{\sum_{t=1}^{N_E} \text{Act}(\text{Router}(\mathbf{x}))[t]} \quad (5)$$

Here,  $\text{argtopK}$  get the index of the top-K routing logits,  $\mathbf{f}_i$  is the fraction of tokens in the batch  $\mathcal{B}$  that are assigned to expert  $i$ , and  $\mathbf{P}_i$  is the average gating weight that the router assigns to expert  $i$  (both estimated over a batch of size  $T$ ). Minimizing  $\mathcal{L}_{\text{aux}}$  therefore penalizes experts that are either under-selected or consistently receive low gating weight, encouraging the router to distribute token assignments and gating weights more evenly across experts. The coefficient  $\alpha_{\text{aux}}$  controls the regularization strength, and the multiplicative factor  $N_E$  normalizes the objective with respect to the number of experts.

#### 3.2 MIXTURE OF NEURON EXPERTS

To explore the sparsity within the activated experts, we decompose the each experts into neuron granular mixture of experts. The output of an expert is formulated as:

$$\mathbf{E}_i(\mathbf{x}) = \mathbf{W}_{\text{down}}^i (\text{SiLU}(\mathbf{W}_{\text{gate}}^i \mathbf{x}) \odot \mathbf{W}_{\text{up}}^i \mathbf{x}) \quad (6)$$

For clarity and compactness, we denote the outputs of the gate projection and the up projection by  $\mathbf{G}$  and  $\mathbf{H}$ , respectively. Concretely,

$$\mathbf{G} = \text{SiLU}(\mathbf{W}_{\text{gate}}^i \mathbf{x}) \in \mathbb{R}^{d_{\text{expert}}}, \quad \mathbf{H} = \mathbf{W}_{\text{up}}^i \mathbf{x} \in \mathbb{R}^{d_{\text{expert}}}. \quad (7)$$

Then Equation (6) can be reformulated as:

$$\mathbf{E}_i(\mathbf{x}) = \mathbf{W}_{\text{down}}^i (\mathbf{G} \odot \mathbf{H}). \quad (8)$$

**Algorithm 1:** Mixture of Neuron Experts (MoNE)

---

```

216 1: Input: Layer input  $\mathbf{x} \in \mathbb{R}^{d_{\text{model}}}$ , number of experts  $n$ , number of selected experts  $K_E$ , number
217   of selected neurons for each expert  $K_N$ , activation function of router  $\text{Act}$ ,
218 2: Output: Layer output  $\mathbf{h} \in \mathbb{R}^{d_{\text{model}}}$ 
219 3:  $\triangleright$  Initialize  $\mathbf{h} = \mathbf{0}$ 
220 4:  $\mathbf{p} = \text{Router}(\mathbf{x}) \in \mathbb{R}^n$  // Calculate the scores for each expert
221 5:  $\mathbf{I}_E = \text{argtopK}(\mathbf{p})$  // Select top- $K_E$  experts
222 6:  $\hat{\mathbf{p}} = \text{Act}(\mathbf{p}[\mathbf{I}_E])$  // Calculate the activated scores
223 7: for each selected expert  $i \in \mathbf{I}_E$  do
224 8:    $\mathbf{G}_i = \text{SiLU}(\mathbf{W}_{\text{gate}}^i \mathbf{x})$  // Calculate the output of down projection
225 9:    $\mathbf{I}_N = \text{argtopK}(\text{Abs}(\mathbf{G}_i))$  // Select top- $K_N$  neurons
226 10:   $\tilde{\mathbf{W}}_{\text{up}}^i = \mathbf{W}_{\text{up}}^i[\mathbf{I}_N, :]$ ,  $\tilde{\mathbf{W}}_{\text{down}}^i = \mathbf{W}_{\text{down}}^i[:, \mathbf{I}_N]$  // Select the weights used for calculation
227 11:   $\tilde{\mathbf{E}}_i(\mathbf{x}) = \tilde{\mathbf{W}}_{\text{down}}^i(\mathbf{G}_i[\mathbf{I}_N] \odot \tilde{\mathbf{W}}_{\text{up}}^i \mathbf{x})$  // Calculate the output of expert  $i$ 
228 12:   $\mathbf{h} = \mathbf{h} + \hat{\mathbf{p}}[i] \cdot \tilde{\mathbf{W}}_{\text{up}}^i \mathbf{x}$  // Sum the layer output
229 13: end for
230 14: return  $\mathbf{h}$ 

```

---

Let  $\mathbf{W}_{\text{down}}^i[:, k] \in \mathbb{R}^{d_{\text{model}}}$  denote the  $k$ -th column of  $\mathbf{W}_{\text{down}}^i$  and  $\mathbf{W}_{\text{up}}^i[k, :] \in \mathbb{R}^{1 \times d_{\text{model}}}$  the  $k$ -th row of  $\mathbf{W}_{\text{up}}^i$ , then we expand the product as follows:

$$\mathbf{E}_i(\mathbf{x}) = \sum_{k=1}^{d_{\text{expert}}} \mathbf{W}_{\text{down}}^i[:, k] (\mathbf{G}[k] \cdot \mathbf{H}[k]) \quad (9)$$

$$= \sum_{k=1}^{d_{\text{expert}}} \mathbf{G}[k] \cdot (\mathbf{W}_{\text{down}}^i[:, k] (\mathbf{W}_{\text{up}}^i[k, :] \mathbf{x})) \quad (10)$$

$$= \sum_{k=1}^{d_{\text{expert}}} \mathbf{G}[k] \cdot \underbrace{(\mathbf{W}_{\text{down}}^i[:, k] \mathbf{W}_{\text{up}}^i[k, :])}_{\mathbf{A}_k} \mathbf{x}. \quad (11)$$

Consequently, we can get the decomposition as follows:

$$\mathbf{E}_i(\mathbf{x}) = \sum_{k=1}^{d_{\text{expert}}} \mathbf{G}[k] \cdot \mathbf{A}_k \mathbf{x}, \quad \text{where } \mathbf{A}_k = \mathbf{W}_{\text{down}}^i[:, k] \mathbf{W}_{\text{up}}^i[k, :] \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}} \quad (12)$$

Eq. equation 12 shows that each expert can be decomposed into a set of neuron granular experts  $\mathbf{A}_k$  that weighted by the neuron level activations  $\mathbf{G}$ . (**In the subsequent articles, we refer to traditional experts as experts and to neuron granular experts as neuron experts.**) Motivated by this perspective, we explore the distribution of  $\mathbf{G}$  in the mainstream MoE models at inference. As shown in Figure 2, the majority of neuron experts receive negligible gate weights: most values of  $\mathbf{G}$  are very small, indicating that a large fraction of neurons inside each expert are inactive during inference. To further quantify the impact of these low-activation neurons, we ablate neuron experts whose gate values fall below a threshold and measure the resulting performance. As shown in Figure 1 we find that retaining approximately the top 30% of neuron experts by gate magnitude is sufficient to preserve the bulk of the original performance, which implies that conventionally trained MoE architectures induce high sparsity in the set of activated parameters at inference.

To address this issue, we propose the *Mixture of Neuron Experts* (MoNE). Algorithm 1 formulates the pipeline of MoNE. Concretely, the router first selects a set of experts as usual; for each selected expert  $i$ , we first calculate the the neuron gating weights  $\mathbf{G}$ , then we use the neuron experts associated with high absolute gating values to calculate the output of each experts. MoNE converts the traditional MoE into a neuron granular MoE via a simple single, per-expert sort-and-select operation. In contrast, an explicit neuron level routing design in traditional MoE would require a substantially larger router and incur heavy cross-expert communication overhead (Lepikhin et al., 2020). MoNE’s selection incurs no additional routing parameters and only accesses the parameters of the expert itself ; because the selected neuron experts communicate within their host expert, the

Table 1: Comparison between traditional MoE and MoNE with the same number of activated experts. MoNE shows comparable results while only use half of the activated parameters in the MoE Layer.

Model	ARC-C	BOOIQ	HELLA	LAMBDA	MNLI	PIQA	RACE	SIQA	WINO	WNLI	AVG.
Traditional MoE 4E/64E	30.55 $\pm$ 1.35	56.94 $\pm$ 0.87	47.78 $\pm$ 0.50	32.70 $\pm$ 0.65	34.39 $\pm$ 0.48	69.53 $\pm$ 0.78	30.33 $\pm$ 0.65	39.87 $\pm$ 1.11	52.80 $\pm$ 1.19	40.85 $\pm$ 1.69	43.57
MoNE w/ Random Selection 4E/64E	27.39 $\pm$ 1.30	56.79 $\pm$ 0.87	37.94 $\pm$ 0.48	27.56 $\pm$ 0.62	35.49 $\pm$ 0.48	65.13 $\pm$ 0.81	30.24 $\pm$ 0.65	38.84 $\pm$ 1.10	49.88 $\pm$ 1.19	43.66 $\pm$ 1.70	42.69
MoNE w/ TopK Selection 4E/64E	31.31 $\pm$ 1.35	53.64 $\pm$ 0.88	48.08 $\pm$ 0.50	34.76 $\pm$ 0.66	33.96 $\pm$ 0.48	70.02 $\pm$ 0.78	31.48 $\pm$ 0.66	39.25 $\pm$ 1.10	51.78 $\pm$ 1.19	47.89 $\pm$ 1.71	<b>44.21</b>

Table 2: Comparison between traditional MoE and MoNE with the same number of activated parameters. For 920M parameter and 2.8M LLMs, MoNE exhibits better downstream performance than MoE models.

Model	ARC-C	BOOIQ	HELLA	LAMBDA	MNLI	PIQA	RACE	SIQA	WINO	WNLI	AVG.
<b>925M Activated 925M</b>											
Dense	33.28 $\pm$ 1.38	58.62 $\pm$ 0.86	52.07 $\pm$ 0.50	37.05 $\pm$ 0.67	33.49 $\pm$ 0.48	71.27 $\pm$ 0.77	30.72 $\pm$ 0.66	40.99 $\pm$ 1.11	54.22 $\pm$ 1.19	52.11 $\pm$ 1.71	46.38
<b>925M Activated 290M</b>											
Traditional MoE 4E/64E	30.55 $\pm$ 1.35	56.94 $\pm$ 0.87	47.78 $\pm$ 0.50	32.70 $\pm$ 0.65	34.39 $\pm$ 0.48	69.53 $\pm$ 0.78	30.33 $\pm$ 0.65	39.87 $\pm$ 1.11	52.80 $\pm$ 1.19	40.85 $\pm$ 1.69	43.57
MoNE w/ o NG-LBL 8E/64E	30.97 $\pm$ 1.35	55.75 $\pm$ 0.87	48.01 $\pm$ 0.50	33.34 $\pm$ 0.66	34.44 $\pm$ 0.48	70.67 $\pm$ 0.78	29.86 $\pm$ 0.65	38.89 $\pm$ 1.10	53.83 $\pm$ 1.19	49.30 $\pm$ 1.72	<b>44.51</b>
MoNE w/ NG-LBL 8E/64E	30.38 $\pm$ 1.34	59.45 $\pm$ 0.86	49.51 $\pm$ 0.50	33.96 $\pm$ 0.66	34.79 $\pm$ 0.48	70.89 $\pm$ 0.77	30.24 $\pm$ 0.65	39.76 $\pm$ 1.11	53.67 $\pm$ 1.19	52.11 $\pm$ 1.71	<b>45.48</b>
<b>925M Activated 310M</b>											
Traditional MoE 6E/64E	33.02 $\pm$ 1.37	54.50 $\pm$ 0.87	49.20 $\pm$ 0.50	34.48 $\pm$ 0.66	35.04 $\pm$ 0.48	71.33 $\pm$ 0.77	30.91 $\pm$ 0.66	40.58 $\pm$ 1.11	53.43 $\pm$ 1.19	36.62 $\pm$ 1.65	43.91
MoNE w/ o NG-LBL 12E/64E	30.97 $\pm$ 1.35	55.99 $\pm$ 0.87	50.07 $\pm$ 0.50	35.73 $\pm$ 0.67	32.35 $\pm$ 0.47	69.97 $\pm$ 0.78	30.81 $\pm$ 0.66	40.43 $\pm$ 1.11	51.78 $\pm$ 1.19	52.11 $\pm$ 1.71	<b>45.02</b>
MoNE w/ NG-LBL 12E/64E	32.17 $\pm$ 1.36	62.11 $\pm$ 0.85	48.65 $\pm$ 0.50	34.58 $\pm$ 0.66	33.89 $\pm$ 0.48	71.44 $\pm$ 0.77	31.00 $\pm$ 0.66	39.20 $\pm$ 1.10	52.88 $\pm$ 1.19	50.70 $\pm$ 1.72	<b>45.75</b>
<b>925M Activated 330M</b>											
Traditional MoE 8E/64E	32.68 $\pm$ 1.37	56.61 $\pm$ 0.87	49.70 $\pm$ 0.50	35.51 $\pm$ 0.67	33.36 $\pm$ 0.48	71.93 $\pm$ 0.77	30.33 $\pm$ 0.65	40.74 $\pm$ 1.11	51.22 $\pm$ 1.19	39.44 $\pm$ 1.68	44.30
MoNE w/ o NG-LBL 16E/64E	31.31 $\pm$ 1.35	56.39 $\pm$ 0.87	50.59 $\pm$ 0.50	35.82 $\pm$ 0.67	35.36 $\pm$ 0.48	71.55 $\pm$ 0.77	31.29 $\pm$ 0.66	41.30 $\pm$ 1.11	52.96 $\pm$ 1.19	52.11 $\pm$ 1.71	<b>45.82</b>
MoNE w/ NG-LBL 16E/64E	30.38 $\pm$ 1.34	60.31 $\pm$ 0.86	49.17 $\pm$ 0.50	34.58 $\pm$ 0.66	34.82 $\pm$ 0.48	70.84 $\pm$ 0.77	30.81 $\pm$ 0.66	40.94 $\pm$ 1.11	51.38 $\pm$ 1.19	50.70 $\pm$ 1.72	<b>45.36</b>
<b>2.81B Activated 0.55B</b>											
Traditional MoE 4E/64E (2.81B)	38.65 $\pm$ 1.42	57.89 $\pm$ 0.87	63.00 $\pm$ 0.48	44.38 $\pm$ 0.69	31.61 $\pm$ 0.47	75.19 $\pm$ 0.74	34.74 $\pm$ 0.68	42.37 $\pm$ 1.12	59.98 $\pm$ 1.17	47.89 $\pm$ 1.71	50.18
MoNE w/ o NG-LBL 8E/64E (2.81B)	39.93 $\pm$ 1.43	61.13 $\pm$ 0.86	63.87 $\pm$ 0.48	46.26 $\pm$ 0.69	38.67 $\pm$ 0.49	76.12 $\pm$ 0.73	34.70 $\pm$ 0.68	42.82 $\pm$ 1.12	59.19 $\pm$ 1.17	43.66 $\pm$ 1.70	<b>50.44</b>
MoNE w/ NG-LBL 8E/64E (2.81B)	37.54 $\pm$ 1.41	62.97 $\pm$ 0.85	63.36 $\pm$ 0.48	46.13 $\pm$ 0.69	36.28 $\pm$ 0.49	76.17 $\pm$ 0.73	34.83 $\pm$ 0.68	42.32 $\pm$ 1.12	59.67 $\pm$ 1.17	57.75 $\pm$ 1.70	<b>51.72</b>

extra communication latency can be negligible. Empirically, pretraining with 50% of the activation parameters by MoNE already matches the performance of a traditional MoE, which demonstrate MoNE can effectively improve the utilization of activated parameters.

### 3.3 NEURON GRANULAR LOAD BALANCE LOSS

Since We decompose each expert into neuron level sub-experts, we further introduce the *neuron granular load balance loss* (NG-LBL). NG-LBL is designed to avoid cases where a subset of neurons are rarely activated, thereby further improving parameter utilization. The formulation of NG-LBL of is similar with  $\mathcal{L}_{\text{aux}}$ , but is applied independently to the neuron experts within each expert.

For expert  $i$ , the fraction of tokens in the batch  $\mathcal{B}$  (with batchsize of  $T$ ) that assigned to neuron  $k$ , and the average gating weights that the gate projection assigned to neuron  $k$  is calculated as follows:

$$\tilde{\mathbf{f}}_{i,k} = \frac{1}{T} \sum_{\mathbf{x} \in \mathcal{B}} \mathbb{1}\{k \in \text{argtopK}(\text{Abs}(\mathbf{G}_i))\}, \quad \tilde{\mathbf{P}}_{i,k} = \frac{1}{T} \sum_{\mathbf{x} \in \mathcal{B}} \frac{\text{Abs}(\mathbf{G}_i[k])}{\sum_{t=1}^{d_{\text{expert}}} \text{Abs}(\mathbf{G}_i[t])}. \quad (13)$$

The whole auxiliary load balance loss used in MoNE is to sum the original auxiliary load balance loss and each expert’s neuron granular load balance loss:

$$\tilde{\mathcal{L}}_{\text{aux}} = \mathcal{L}_{\text{aux}} + \sum_{i=1}^{N_E} \mathcal{L}_{\text{NG-LBL}}^i, \quad \text{where} \quad \mathcal{L}_{\text{NG-LBL}}^i = \alpha_{\text{NG-LBL}} \cdot d_{\text{model}} \cdot \sum_{k=1}^{d_{\text{model}}} \tilde{\mathbf{f}}_{i,k} \cdot \tilde{\mathbf{P}}_{i,k} \quad (14)$$

where the  $\alpha_{\text{NG-LBL}}$  is the coefficient to controls the regularization strength of NG-LBL.

## 4 EXPERIMENT

### 4.1 EXPERIMENTAL SETUP

**Model Architectures.** As shown in Table 6, we implement models with total parameter counts of 925M and 2.81B. The MoE layers for both model scales contained 64 experts. Follow by Dai et al.

Table 3: Architecture exploration on different numbers of selected neurons  $K_N$ . MoNE exhibits better downstream performance when the selected rate is  $1/4$ .

Model	ARC-C	BOOIQ	HELLA	LAMBDA	MNLI	PIQA	RACE	SIQA	WINO	WNLI	AVG.
<i>2.81B Activated 0.55B</i>											
Traditional MoE 4E/64E	38.65 $\pm$ 1.42	57.89 $\pm$ 0.87	63.00 $\pm$ 0.48	44.38 $\pm$ 0.69	31.61 $\pm$ 0.47	75.19 $\pm$ 0.74	34.74 $\pm$ 0.68	42.37 $\pm$ 1.12	59.98 $\pm$ 1.17	47.89 $\pm$ 1.71	49.57
$K_N = 1/2 \cdot d_{\text{model}}$											
MoNE w/o NG-LBL 6E/64E	41.04 $\pm$ 1.44	57.34 $\pm$ 0.87	63.50 $\pm$ 0.48	46.21 $\pm$ 0.69	36.98 $\pm$ 0.49	75.68 $\pm$ 0.73	35.12 $\pm$ 0.68	43.35 $\pm$ 1.12	59.76 $\pm$ 1.17	45.07 $\pm$ 1.71	50.41
MoNE w/ NG-LBL 6E/64E	39.59 $\pm$ 1.43	61.96 $\pm$ 0.85	63.10 $\pm$ 0.48	44.89 $\pm$ 0.69	38.41 $\pm$ 0.49	75.30 $\pm$ 0.73	35.22 $\pm$ 0.68	42.17 $\pm$ 1.12	59.12 $\pm$ 1.17	45.07 $\pm$ 1.71	<b>50.48</b>
$K_N = 1/4 \cdot d_{\text{model}}$											
MoNE w/o NG-LBL 8E/64E	39.93 $\pm$ 1.43	61.13 $\pm$ 0.86	63.87 $\pm$ 0.48	46.26 $\pm$ 0.69	38.67 $\pm$ 0.49	76.12 $\pm$ 0.73	34.70 $\pm$ 0.68	42.82 $\pm$ 1.12	59.19 $\pm$ 1.17	43.66 $\pm$ 1.70	50.64
MoNE w/ NG-LBL 8E/64E	37.54 $\pm$ 1.41	62.97 $\pm$ 0.85	63.36 $\pm$ 0.48	46.13 $\pm$ 0.69	36.28 $\pm$ 0.49	76.17 $\pm$ 0.73	34.83 $\pm$ 0.68	42.32 $\pm$ 1.12	59.67 $\pm$ 1.17	57.75 $\pm$ 1.70	<b>51.70</b>
$K_N = 1/10 \cdot d_{\text{model}}$											
MoNE w/o NG-LBL 10E/64E	40.27 $\pm$ 1.43	60.86 $\pm$ 0.86	63.89 $\pm$ 0.48	46.09 $\pm$ 0.69	31.59 $\pm$ 0.47	75.84 $\pm$ 0.73	34.55 $\pm$ 0.68	43.09 $\pm$ 1.12	58.01 $\pm$ 1.17	45.07 $\pm$ 1.71	49.93
MoNE w/ NG-LBL 10E/64E	39.85 $\pm$ 1.43	58.10 $\pm$ 0.87	61.99 $\pm$ 0.48	44.21 $\pm$ 0.69	32.07 $\pm$ 0.47	75.03 $\pm$ 0.74	35.12 $\pm$ 0.68	42.12 $\pm$ 1.12	59.27 $\pm$ 1.17	57.75 $\pm$ 1.70	<b>50.55</b>

Table 4: The performance of MoNE when applying different activation functions. MoNE exhibits better downstream performance when applying SiLU and Sigmoid.

Model	ARC-C	BOOIQ	HELLA	LAMBDA	MNLI	PIQA	RACE	SIQA	WINO	WNLI	AVG.
<i>SiLU</i>											
Traditional MoE 4E/64E	30.55 $\pm$ 1.35	56.94 $\pm$ 0.87	47.78 $\pm$ 0.50	32.70 $\pm$ 0.65	34.39 $\pm$ 0.48	69.53 $\pm$ 0.78	30.33 $\pm$ 0.65	39.87 $\pm$ 1.11	52.80 $\pm$ 1.19	40.85 $\pm$ 1.69	43.38
MoNE w/o NG-LBL 8E/64E	30.97 $\pm$ 1.35	55.75 $\pm$ 0.87	48.01 $\pm$ 0.50	33.34 $\pm$ 0.66	34.44 $\pm$ 0.48	70.67 $\pm$ 0.78	29.86 $\pm$ 0.65	38.89 $\pm$ 1.10	53.83 $\pm$ 1.19	49.30 $\pm$ 1.72	44.67
MoNE w/ NG-LBL 8E/64E	30.38 $\pm$ 1.34	59.45 $\pm$ 0.86	49.51 $\pm$ 0.50	33.96 $\pm$ 0.66	34.79 $\pm$ 0.48	70.89 $\pm$ 0.77	30.24 $\pm$ 0.65	39.76 $\pm$ 1.11	53.67 $\pm$ 1.19	52.11 $\pm$ 1.71	<b>45.48</b>
<i>Sigmoid</i>											
Traditional MoE 4E/64E	30.38 $\pm$ 1.35	54.79 $\pm$ 0.87	45.78 $\pm$ 0.50	33.28 $\pm$ 0.65	34.44 $\pm$ 0.48	69.97 $\pm$ 0.78	30.24 $\pm$ 0.65	38.87 $\pm$ 1.11	53.10 $\pm$ 1.19	43.66 $\pm$ 1.70	43.15
MoNE w/o NG-LBL 8E/64E	31.40 $\pm$ 1.36	49.45 $\pm$ 0.88	48.93 $\pm$ 0.50	34.76 $\pm$ 0.66	35.00 $\pm$ 0.48	71.33 $\pm$ 0.77	32.73 $\pm$ 0.67	39.51 $\pm$ 1.11	49.57 $\pm$ 1.19	50.70 $\pm$ 1.72	44.34
MoNE w/ NG-LBL 8E/64E (Sigmoid)	30.72 $\pm$ 1.35	59.79 $\pm$ 0.86	47.51 $\pm$ 0.50	33.75 $\pm$ 0.66	34.80 $\pm$ 0.48	70.18 $\pm$ 0.78	32.82 $\pm$ 0.67	40.17 $\pm$ 1.11	51.38 $\pm$ 1.19	53.52 $\pm$ 1.71	<b>45.46</b>
<i>Softmax</i>											
Traditional MoE 4E/64E	27.93 $\pm$ 1.31	48.20 $\pm$ 0.88	36.67 $\pm$ 0.48	29.82 $\pm$ 0.64	33.28 $\pm$ 0.48	66.51 $\pm$ 0.80	29.96 $\pm$ 0.65	37.84 $\pm$ 1.10	49.28 $\pm$ 1.19	40.66 $\pm$ 1.69	40.02
MoNE w/o NG-LBL 8E/64E	28.24 $\pm$ 1.31	47.52 $\pm$ 0.88	38.30 $\pm$ 0.49	30.08 $\pm$ 0.64	35.00 $\pm$ 0.48	66.43 $\pm$ 0.80	30.33 $\pm$ 0.65	39.25 $\pm$ 1.10	51.54 $\pm$ 1.19	42.25 $\pm$ 1.70	40.96
MoNE w/ NG-LBL 8E/64E	28.58 $\pm$ 1.32	48.53 $\pm$ 0.88	44.96 $\pm$ 0.50	32.52 $\pm$ 0.65	35.16 $\pm$ 0.48	69.91 $\pm$ 0.78	30.43 $\pm$ 0.66	38.89 $\pm$ 1.10	50.59 $\pm$ 1.19	46.48 $\pm$ 1.71	<b>41.81</b>

(2024), we use one shared expert. The MoE layers in both model scales comprise 64 experts. For the smaller model, we trained traditional MoE with 4, 6 and 8 experts-per-token, which correspond to activated parameters of 290M, 310M and 330M, respectively; For the larger model, we trained a traditional MoE with 4 experts-per-token, corresponding to 0.55B activated parameters. In the main experimental suite, each traditional MoE configuration was paired with a corresponding MoNE: the  $K_N$  is set to  $d_{\text{model}}/4$  and  $K_E$  is set two times of corresponding traditional MoE. so that the number of activated parameters is equal between the traditional MoE and MoNE. Please refer Appendix A.1 for detailed training hyper-parameters.

**Data & Tokenizer.** We trained the 925M parameter model on a 50B-token subset of the NeMaTron-CC dataset (Su et al., 2024) and the 2.81B parameter model on a 100B-token subset of the NeMaTron-CC dataset. All text was tokenized with the LLaMA-3-8B tokenizer (Dubey et al., 2024).

**Hyper-Parameters.** The hyper-parameters are selected based on the common practice for dense language models. We replace all FFN layers with MoE layer in the transformer. Please refer Appendix A.1 for detailed training hyper-parameters.

**Benchmarks.** We use the lm-evaluation-harness (Gao et al., 2024) for evaluation. The benchmarks used include ARC-C (Clark et al., 2018), BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), LAMBADA (Paperno et al., 2016), MNLI (Williams et al., 2017), PIQA (Bisk et al., 2020), RACE (Lai et al., 2017), SIQA (Sap et al., 2019), WinoGrande (Sakaguchi et al., 2021), WNLI (Levesque et al., 2012). For all these benchmarks, we report the zero-shot accuracy.

## 4.2 MAIN RESULTS

**Prior Experiment** We first compare traditional MoE and MoNE with the same number of experts per token  $N_E$ . For MoNE, we set the number of used neurons  $K_N$  to  $d_{\text{model}}/4$ , which reduces the number of MoNE’s activated parameters in the moe layer to approximately half of the traditional MoE’s. As an additional baseline, we select a random subset of size  $K_N = d_{\text{model}}/4$  from the neuron experts instead of the  $\text{TOP}K$  strategy to pretrain the model. Table 1 shows that MoNE attains performance comparable to the traditional MoE while using 50% of the activated parameters in the

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

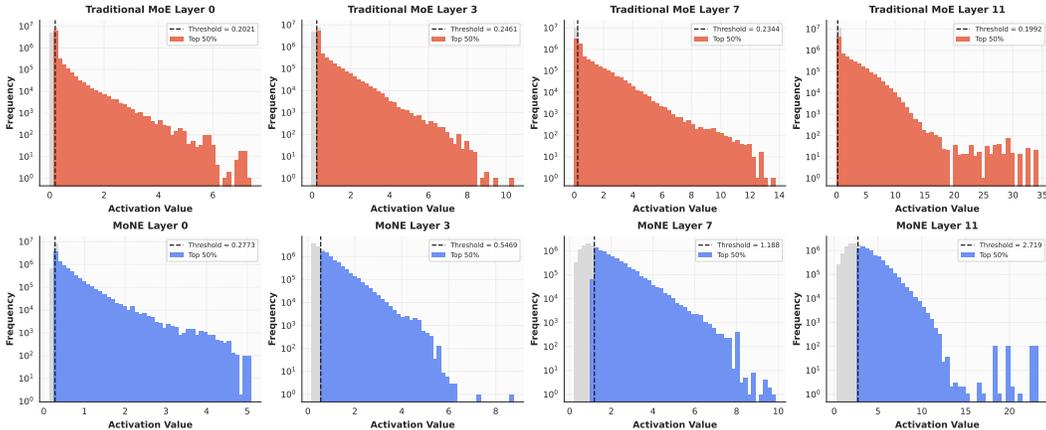


Figure 4: The comparison of the activation value  $G$  for the neuron experts between traditional MoE and MoNE. MoNE effectively increase the activation weight compared with traditional MoE.

MoE layer, whereas the random subset baseline suffers a large drop in performance. These findings indicate that MoNE’s selection mechanism effectively improve the utilization of the activated parameters.

**Comparisons to Traditional MoE of Equivalent Activated Parameters** We further compared MoNE with Traditional MoE with equivalent activated parameters. As shown in the Table 2, MoNE consistently improves over the traditional MoE, and the improvement grows as the number of activated parameters increases. Concretely, when activating 290M parameters MoNE yields an improvement of approximately 1% relative to the traditional MoE, and this gain increases to about 2% at 330M activated parameters. Also, with 330M activated parameters, MoNE attains performance that is comparable to a dense model with 925M parameters. For the 3B models MoNE improves upon the traditional MoE by roughly 1.1%. These results indicate that MoNE is a promising approach for training MoE-like models.

### 4.3 FURTHER ANALYSIS

**The Effectiveness of NG-LBL** We investigate the effect of NG-LBL on MoNE in Table 2 and Table 3. Empirically, NG-LBL consistently improves MoNE’s performance and increases its parameter efficiency: on the 1B-parameter model, NG-LBL yields an improvement about 1.0%, while on a 3B-parameter model the gain is about 1.4%. Figure 5 shows that NG-LBL substantially accelerates the decline of training loss, which demonstrates the effectiveness of NG-LBL. To better understand how NG-LBL helps, we examine load balancing among neuron experts. As shown in Figure 6, neuron experts achieve better load balance with NG-LBL. This indicates that the balancing at the neuron level can effectively increase the ability of experts

**The Analysis of the Activation Weight** Figure 4 compares the activation weight of traditional MoE and MoNE. MoNE effectively increases neuron activation weight: the median value of activation weight in the first layer rises from 0.20 to 0.28 and continues to grow with depth. The median value increases from 0.20 to 2.70 in the final layer. In addition, the activation weight’s distribution produced by MoNE is noticeably more uniform, indicating a more balanced utilization of neuron experts. These results demonstrate that MoNE both increases and homogenizes the use of neuron experts, thereby reducing the sparsity of activated parameters and improving the utilization of the activated parameters.

**The Effect of the Neuron Expert Activation Ratio** Figure 1 indicates that the within-expert activation rate has a meaningful effect. Keeping the number of activated parameters fixed, we pretrained models with different  $K_N$ ; Table 3 shows that all settings improve over the baseline, with the best performance at a ratio of  $1/4$ . Therefore, we recommend using  $K_N = 1/4 \cdot d_{\text{model}}$ . The result aligns with the result in Figure 1: model performance remains stable until approximately 70% of neurons

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

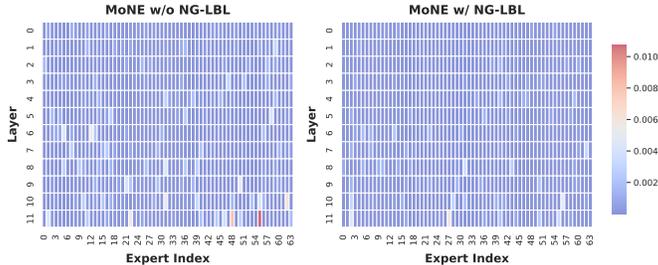
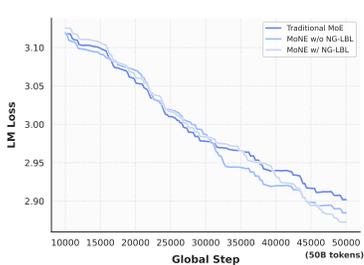


Figure 5: Pre-training loss between traditional MoE and MoNE Figure 6: The visualization of load balance for different layers and experts, the value visualized is the variance of  $\tilde{f}_{i,k}$ .

are removed, which implies that each activated expert only needs 30% of its neurons to process an input.

**The Effect of Different Activation Function inside the Expert** We further explored three internal activation functions for neuron experts in MoNE—Sigmoid, SiLU, and Softmax. Empirically, Sigmoid and SiLU produce consistently performance than Softmax. While experts contain large amounts of neurons, Softmax concentrates probability distribution on a small subset while assigning near-zero weights to most neurons, thereby reducing effective parameter usage. Using NG-LBL mitigates this concentration by encouraging more uniform neuron activation, but Softmax still underperforms the baseline in our experiments. Accordingly, we recommend use Sigmoid or SiLU as the default internal activation for MoNE architectures.

**The Efficiency of MoNE** We further investigate the efficiency of MoNE. As shown in Table 5, with the same number of activated parameters, MoNE and traditional MoE require comparable GPU memory and show nearly identical throughput. Crucially, MoNE achieves neuron granular expert selection without enlarging the router or increasing communication latency overhead. Hence, MoNE provides a practical approach to neuron granular expert computation while preserving computation efficiency.

## 5 CONCLUSION

In this work, we demonstrate that the parameters activated by Mixture-of-Experts (MoE) layers is also highly sparse. By decomposing each expert into neuron granular subexperts, we find that many neuron experts receive very small activation weights. The result motivate us to improve the utilization of activated parameters by only use the neuron experts with high activation weights. Therefore We propose Mixture of Neuron Experts (MoNE), a simple and practical modification of traditional MoE that operates at neuron granularity: by decomposing experts into neuron granular subexperts and applying a simple sorting operation to the gate-projection outputs prior to expert computation, MoNE converts a traditional MoE into a neuron granular MoE. Furthermore, we propose to apply the neuron granular load-balance loss on the neuron experts to encourage more uniform neuron utilization. MoNE requires no additional model parameters and incurs only a negligible computational overhead relative to traditional MoE. Empirically, MoNE matches baseline performance while activating only half of the parameters in the MoE layer and achieves consistent improvements when compared at equal numbers of activated parameters. Expert granularity is an important focus of current MoE development, while traditional MoE faces problems such as large routers and large

Table 5: Throughput and memory usage comparison among traditional MoE and MoNE. The model with 925M parameters and 290M activation parameters is used for calculation. Auxiliary losses do not impact efficiency.

	Traditional MoE	MoNE
<b>Configuration</b>		
Batch size	8	8
Input length	1024	1024
New tokens	128	128
<b>Throughput &amp; Memory</b>		
Tokens/sec	1340.63	1338.49
Memory Peak Reserved	7.8GB	7.8GB

486 communication delays when expert partitioning is overly fine granularity,. We believe MoNE is a  
487 practical step toward more efficient and scalable MoE-like architectures.  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

540 ETHICS STATEMENT

541

542 This paper presents work whose goal is to advance the field of large language model. There are  
543 many potential consequences of our work, none of which we feel must be specifically highlighted  
544 here.

545

546 REPRODUCIBILITY STATEMENT

547

548 The details of datasets, model architectures and hyper-parameters are described in Section 4.1 and  
549 Appendix A.1.

550

551 REFERENCES

552

553 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
554 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
555 report. *arXiv preprint arXiv:2303.08774*, 2023.

556

557 Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K  
558 Arora, Yu Bai, Bowen Baker, Haiming Bao, et al. gpt-oss-120b & gpt-oss-20b model card. *arXiv  
559 preprint arXiv:2508.10925*, 2025.

560 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge,  
561 Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

562

563 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical com-  
564 monsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*,  
565 volume 34, pp. 7432–7439, 2020.

566 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
567 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
568 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

569

570 Weilin Cai, Juyong Jiang, Fan Wang, Jing Tang, Sunghun Kim, and Jiayi Huang. A survey on  
571 mixture of experts in large language models. *IEEE Transactions on Knowledge and Data Engi-  
572 neering*, 2025.

573 Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina  
574 Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint  
575 arXiv:1905.10044*, 2019.

576 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
577 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
578 *arXiv preprint arXiv:1803.05457*, 2018.

579

580 Damai Dai, Chengqi Deng, Chenggang Zhao, RX Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding  
581 Zeng, Xingkai Yu, Yu Wu, et al. Deepseekmoe: Towards ultimate expert specialization in mixture-  
582 of-experts language models. *arXiv preprint arXiv:2401.06066*, 2024.

583 Yann N Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated  
584 convolutional networks. In *International conference on machine learning*, pp. 933–941. PMLR,  
585 2017.

586

587 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep  
588 bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of  
589 the North American chapter of the association for computational linguistics: human language  
590 technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.

591 Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim  
592 Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of language  
593 models with mixture-of-experts. In *International conference on machine learning*, pp. 5547–  
5569. PMLR, 2022.

- 594 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
595 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
596 *arXiv e-prints*, pp. arXiv-2407, 2024.
- 597  
598 William Fedus, Jeff Dean, and Barret Zoph. A review of sparse expert models in deep learning.  
599 *arXiv preprint arXiv:2209.01667*, 2022a.
- 600  
601 William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter  
602 models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39,  
603 2022b.
- 604  
605 Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural  
606 networks. *arXiv preprint arXiv:1803.03635*, 2018.
- 607  
608 Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M Roy, and Michael Carbin. Stabilizing the  
609 lottery ticket hypothesis. *arXiv preprint arXiv:1903.01611*, 2019.
- 610  
611 Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in  
612 one-shot. In *International conference on machine learning*, pp. 10323–10337. PMLR, 2023.
- 613  
614 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster,  
615 Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muen-  
616 nighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang  
617 Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model  
618 evaluation harness, 07 2024. URL <https://zenodo.org/records/12608602>.
- 619  
620 Xinyang Geng and Hao Liu. Openllama: An open reproduction of llama, May 2023. URL [https://github.com/openlm-research/open\\_llama](https://github.com/openlm-research/open_llama).
- 621  
622 Xu Owen He. Mixture of a million experts. *arXiv preprint arXiv:2407.04153*, 2024.
- 623  
624 Zihao Huang, Qiyang Min, Hongzhi Huang, Defa Zhu, Yutao Zeng, Ran Guo, and Xun Zhou. Ultra-  
625 sparse memory network. *arXiv preprint arXiv:2411.12364*, 2024.
- 626  
627 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bam-  
628 ford, Devendra Singh Chaitan, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.  
629 Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- 630  
631 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,  
632 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
633 models. *arXiv preprint arXiv:2001.08361*, 2020.
- 634  
635 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,  
636 2014.
- 637  
638 Jakub Krajewski, Jan Ludziejewski, Kamil Adamczewski, Maciej Pióro, Michał Krutul, Szymon  
639 Antoniuk, Kamil Ciebiera, Krystian Król, Tomasz Odrzygóźdź, Piotr Sankowski, et al. Scaling  
640 laws for fine-grained mixture of experts. *arXiv preprint arXiv:2402.07871*, 2024.
- 641  
642 Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading  
643 comprehension dataset from examinations. *arXiv preprint arXiv:1704.04683*, 2017.
- 644  
645 Guillaume Lample, Alexandre Sablayrolles, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé  
646 Jégou. Large memory layers with product keys. *Advances in Neural Information Processing  
647 Systems*, 32, 2019.
- 648  
649 Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang,  
650 Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional  
651 computation and automatic sharding. *arXiv preprint arXiv:2006.16668*, 2020.
- 652  
653 Hector J Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. *KR*,  
654 2012(13th):3, 2012.

- 648 Zichong Li, Chen Liang, Zixuan Zhang, Ilgee Hong, Young Jin Kim, Weizhu Chen, and Tuo Zhao.  
649 Slimmoe: Structured compression of large moe models via expert slimming and distillation. *arXiv*  
650 *preprint arXiv:2506.18349*, 2025.
- 651 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,  
652 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*  
653 *arXiv:2412.19437*, 2024.
- 654 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*  
655 *arXiv:1711.05101*, 2017.
- 656 Saeed Masoudnia and Reza Ebrahimpour. Mixture of experts: a literature survey. *Artificial Intelli-*  
657 *gence Review*, 42(2):275–293, 2014.
- 658 Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi,  
659 Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset:  
660 Word prediction requiring a broad discourse context. *arXiv preprint arXiv:1606.06031*, 2016.
- 661 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
662 Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text  
663 transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- 664 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adver-  
665 sarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- 666 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialliqa: Common-  
667 sense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- 668 Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman  
669 Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. BLOOM: A 176b-  
670 parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.
- 671 N Shazeer, A Mirhoseini, K Maziarz, A Davis, Q Le, G Hinton, and J Dean. The sparsely-gated  
672 mixture-of-experts layer. *Outrageously large neural networks*, 2, 2017.
- 673 Dan Su, Kezhi Kong, Ying Lin, Joseph Jennings, Brandon Norrick, Markus Kliegl, Mostofa Patwary,  
674 Mohammad Shoeybi, and Bryan Catanzaro. Nemotron-cc: Transforming common crawl into a  
675 refined long-horizon pretraining dataset. *arXiv preprint arXiv:2412.02595*, 2024.
- 676 Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen,  
677 Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv*  
678 *preprint arXiv:2507.20534*, 2025.
- 679 Qwen Team. Qwen3 technical report, 2025.
- 680 Changxin Tian, Kunlong Chen, Jia Liu, Ziqi Liu, Zhiqiang Zhang, and Jun Zhou. Towards  
681 greater leverage: Scaling laws for efficient mixture-of-experts language models. *arXiv preprint*  
682 *arXiv:2507.17702*, 2025.
- 683 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
684 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
685 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 686 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
687 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-  
688 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 689 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yo-  
690 gatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language  
691 models. *arXiv preprint arXiv:2206.07682*, 2022.
- 692 Adina Williams, Nikita Nangia, and Samuel R Bowman. A broad-coverage challenge corpus for  
693 sentence understanding through inference. *arXiv preprint arXiv:1704.05426*, 2017.

702 Haoyuan Wu, Haoxing Chen, Xiaodong Chen, Zhanchao Zhou, Tiejuan Chen, Yihong Zhuang,  
703 Guoshan Lu, Zenan Huang, Junbo Zhao, Lin Liu, et al. Grove moe: Towards efficient and superior  
704 moe llms with adjugate experts. *arXiv preprint arXiv:2508.07785*, 2025.  
705  
706 Fuzhao Xue, Zian Zheng, Yao Fu, Jinjie Ni, Zangwei Zheng, Wangchunshu Zhou, and Yang  
707 You. Openmoe: An early effort on open mixture-of-experts language models. *arXiv preprint*  
708 *arXiv:2402.01739*, 2024.  
709  
710 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-  
711 chine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.  
712  
713 Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small  
714 language model. *arXiv preprint arXiv:2401.02385*, 2024.  
715  
716 Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christo-  
717 pher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. OPT: Open pre-trained transformer  
718 language models. *arXiv preprint arXiv:2205.01068*, 2022.  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A APPENDIX

### A.1 EXPERIMENTAL DETAILS

#### A.1.1 MODEL ARCHITECTURES.

We list the model configuration in Table 6. Here we verified the corresponding MoE and MoNE has the same number of activated parameters. Suppose the number of parameters for `gate projection`, `up projection` and `down projection` is  $N$ . For a MoE layer with 4 experts activated, the number of activated parameters is  $4 \cdot 3 \cdot N$ . For a MoNE layer with 6 experts activated and  $N_k/d_{\text{model}}$  is  $1/2$ , the number of activated parameters of an expert can be calculated as  $N + 1/2 \cdot 2 \cdot N$ , where the parameters of `gate projection` are all activated, and the parameters of `up projection` and `down projection` only activated  $1/2$ . Then total activated parameters can be calculated as  $6 \cdot (N + 1/2 \cdot 2 \cdot N) = 12N$ . Accordingly, for a MoNE layer with 8 experts activated and  $N_k/d_{\text{model}}$  is  $1/4$ , the number of activated parameters is  $8 \cdot (N + 1/4 \cdot 2 \cdot N) = 12N$ , and for a MoNE layer with 10 experts activated and  $N_k/d_{\text{model}}$  is  $1/10$ , the number of activated parameters is  $10 \cdot (N + 1/10 \cdot 2 \cdot N) = 12N$ . Therefore the corresponding MoE and MoNE has the same number of activated parameters.

Table 6: **Sizes and architectures of MoNE and traditional MoE models.** “290M/925M” represents an architecture of an approximately 925M parameter, with 290M activated per token during inference.

Methods	# Layers	# Hidden Size	# Intermediate Size	# Heads	# Head Dim	# The Number of FFN Experts	# The Number of Experts per Token	# $N_k/d_{\text{model}}$
Traditional MoE 290M/925M	12	768	368	16	48	64	4	-
Traditional MoE 310M/925M	12	768	368	16	48	64	6	-
Traditional MoE 330M/925M	12	768	368	16	48	64	8	-
Traditional MoE 0.55B/2.81B	24	1024	512	16	96	64	4	-
MoNE 290M/925M	12	768	368	16	48	64	8	1/4
MoNE 310M/925M	12	768	368	16	48	64	12	1/4
MoNE 330M/925M	12	768	368	16	48	64	16	1/4
MoNE 0.55B/2.81B 6E	24	1024	512	16	96	64	6	1/2
MoNE 0.55B/2.81B 8E	24	1024	512	16	96	64	8	1/4
MoNE 0.55B/2.81B 10E	24	1024	512	16	96	64	10	1/10

#### A.1.2 HYPER-PARAMETERS.

The hyperparameters are selected based on the common practice for dense transformer language models (Zhang et al., 2024; Geng & Liu, 2023; Touvron et al., 2023b; Xue et al., 2024). The key training hyperparameters used in our experiments are as follows: batch size (tokens) = 1M; auxiliary load-balance weight  $\alpha_{\text{aux}} = 0.001$ ; neuron-granular load-balance weight  $\alpha_{\text{NG}} = 0.001$ ; optimizer = FusedAdam (Kingma, 2014); learning rate =  $5e - 4$ ; router scoring activation function = softmax; weight decay (Loshchilov & Hutter, 2017) = 0.1; model maximum sequence length = 2k. These settings were kept fixed across the reported pretraining runs and ablations unless stated otherwise.

#### A.1.3 CALCULATE RESOURCES AND ENVIRONMENT

We use deepspeed as the training framework. For the 925M model, We conduct training on a cluster with 4 nodes and 32 A100 GPUs. For the 2.81B model, We conduct training on a cluster with 16 nodes and 128 A100 GPUs.

## A.2 ADDITIONAL EXPERIMENT

## A.2.1 MORE RESULTS ON THE ACTIVATION VALUE FOR THE NEURON EXPERTS.

We visualize additional neuron granular activation values for Qwen3-30B-A3B and DeepSeek-V2-Lite. As shown in Figure 7 and Figure 8, the vast majority of neuron experts receive negligible gate weights: most entries of  $G$  are close to zero, indicating that a large fraction of neurons within each expert remain effectively inactive during inference.

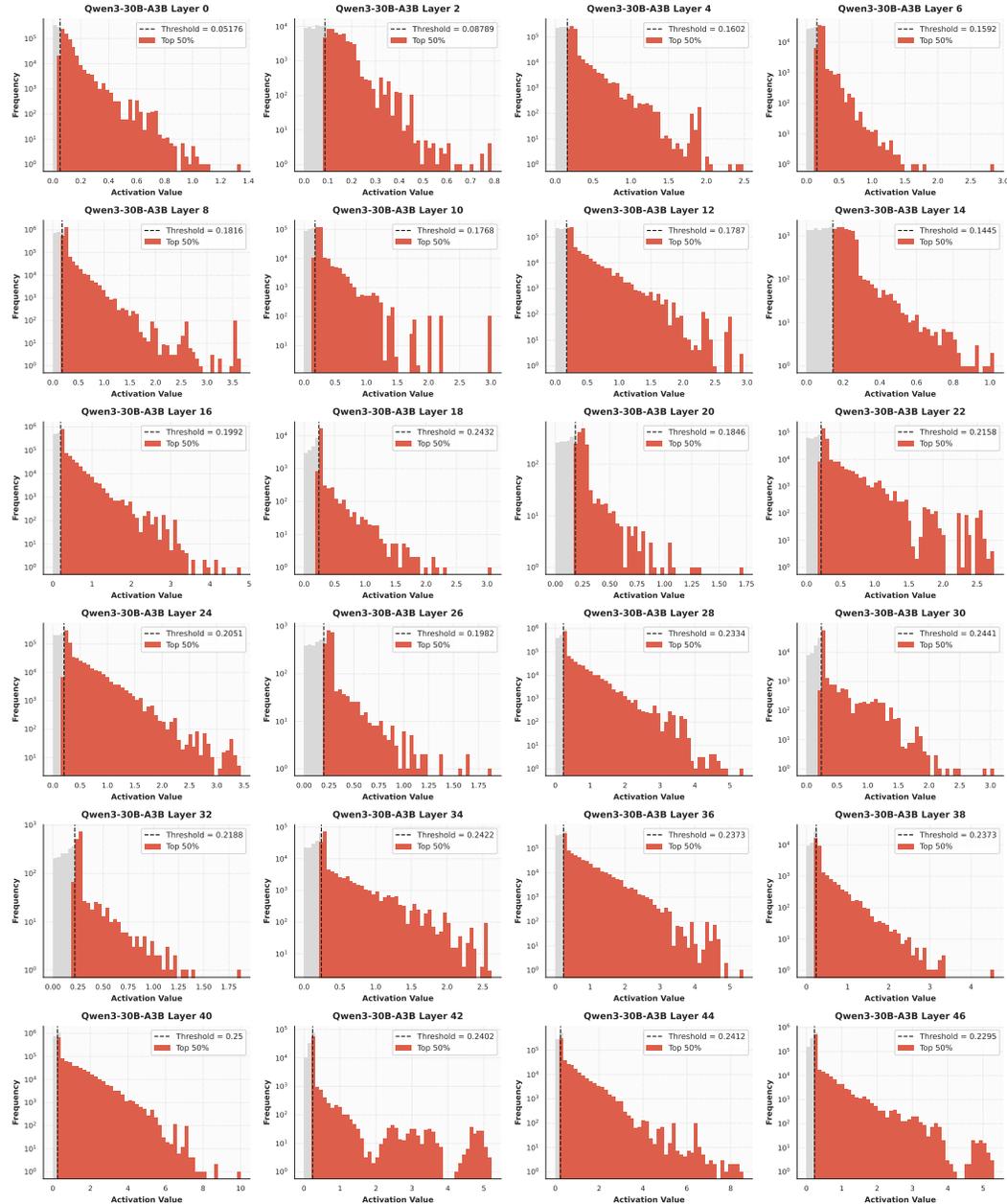


Figure 7: The activation value for the neuron experts on Qwen3-30B-A3

We further compares the activation weight of traditional MoE and MoNE. As shown in Figure 9 and Figure 10 MoNE effectively increases neuron activation weight: the median value of activation weight in the first layer rises from 0.20 to 0.28 and continues to grow with depth. The median

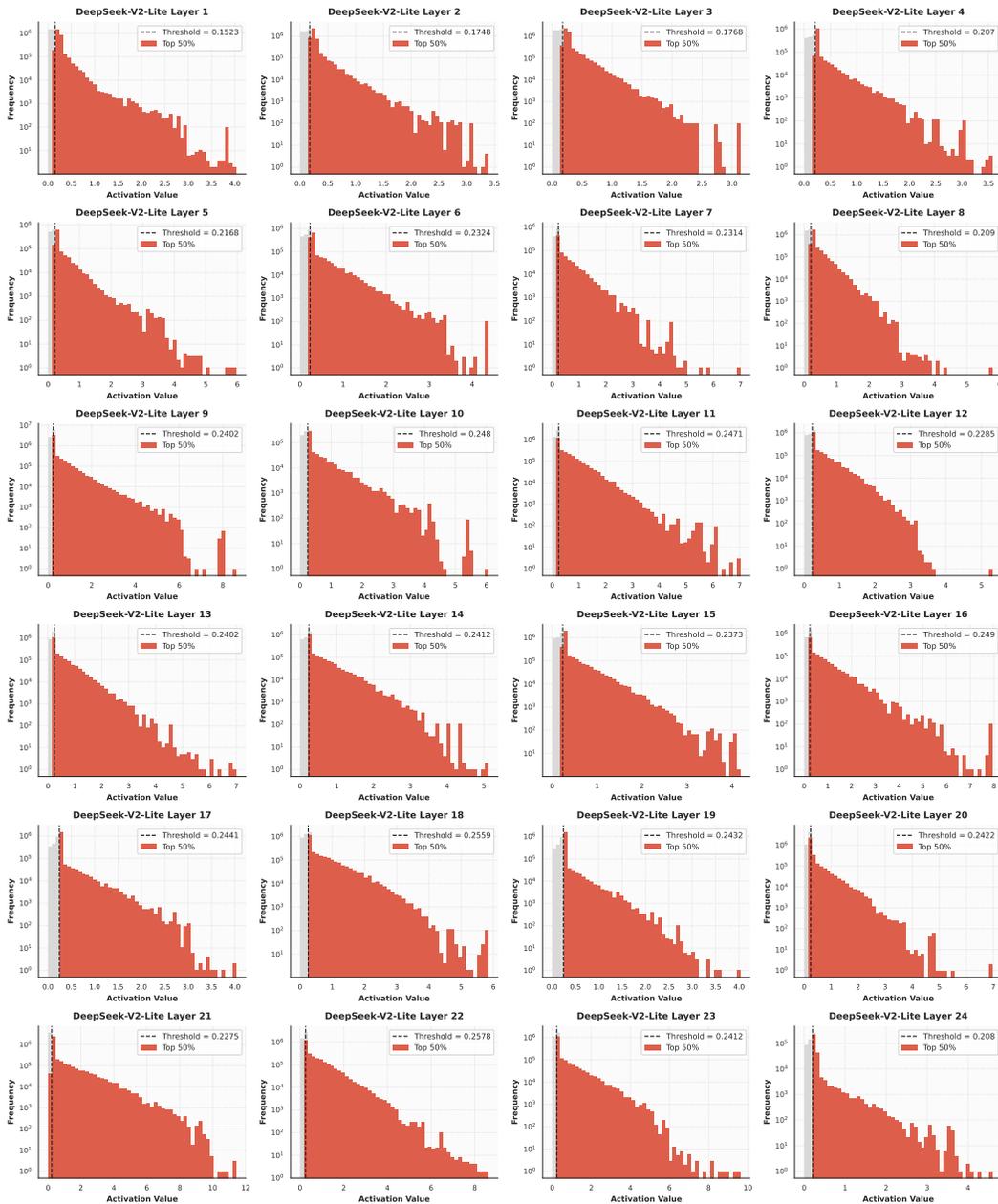


Figure 8: The activation value for the neuron experts on DeepSeek-V2-Lite

value increases from 0.20 to 2.70 in the final layer. In addition, the activation weight’s distribution produced by MoNE is noticeably more uniform, indicating a more balanced utilization of neuron experts. These results demonstrate that MoNE both increases and homogenizes the use of neuron experts, thereby reducing the sparsity of activated parameters and improving the utilization of the activated parameters.

### A.2.2 THE COMPARISON OF TRAINING LOSS FOR TRADITIONAL MOE AND MoNE

As shown in Figure 11, MoNE exhibits more effective expert learning compared with traditional MoE, as evidenced by lower loss values.

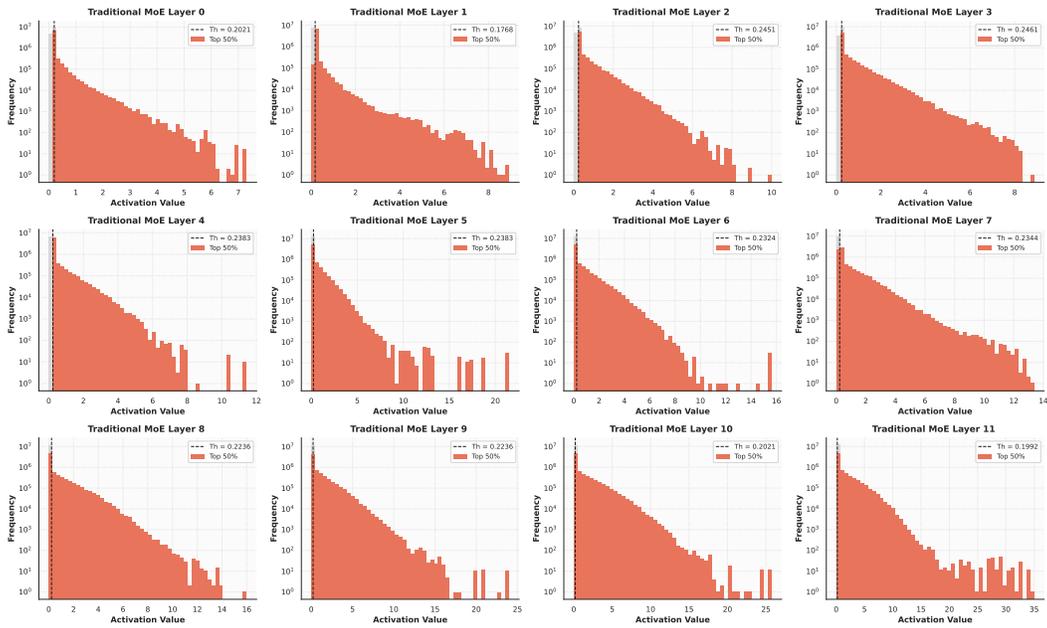


Figure 9: The activation value for the neuron experts on Traditional MoE

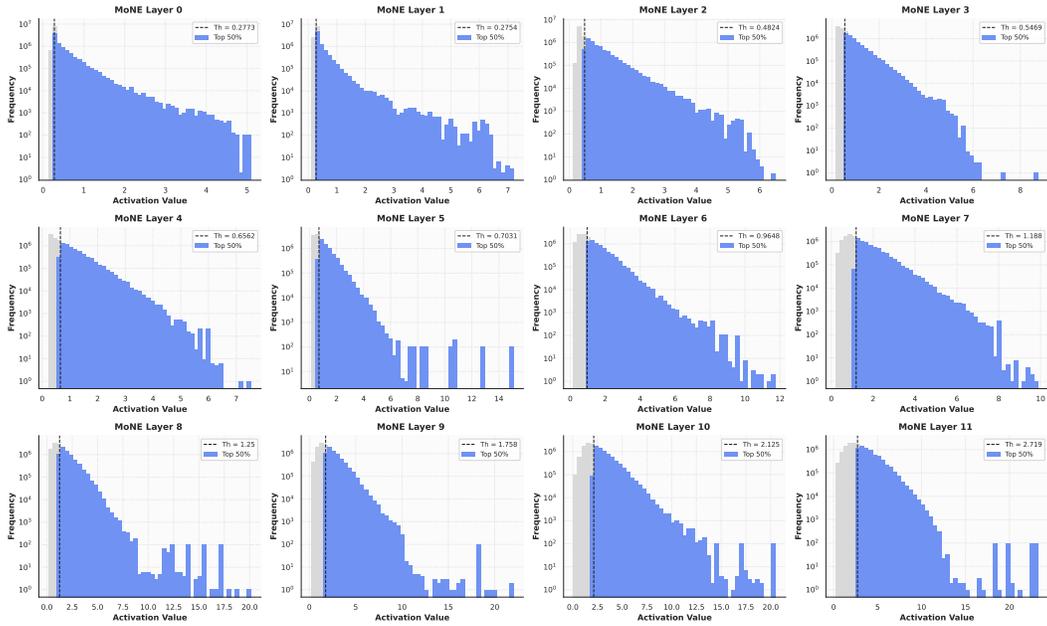


Figure 10: The activation value for the neuron experts on MoNE

### A.2.3 THE EFFECT OF $\mathcal{L}_{\text{AUX}}$ ON MoNE

We further investigate the influence of an auxiliary load-balance loss  $\mathcal{L}_{\text{aux}}$  on MoNE. Our experimental results show that  $\mathcal{L}_{\text{aux}}$  significantly affects MoNE’s performance, suggesting that the balance across experts is important for MoNE to realize further gains in parameter utilization.

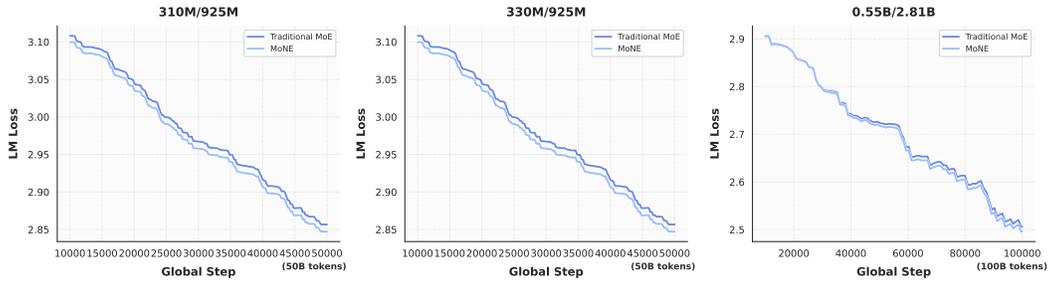


Figure 11: Pre-training loss between tradition MoE and MoNE.

Table 7: The ablation study on auxiliary load-balance loss  $\mathcal{L}_{aux}$ 

Model	ARC-C	BOOIQ	HELLA	LAMBDA	MNLI	PIQA	RACE	SIQA	WINO	WNLI	AVG.
<i>925M Activated 290M</i>											
Traditional MoE w/o $\mathcal{L}_{aux}$ 4E/64E	28.33	46.85	45.63	33.44	34.58	68.88	30.72	38.69	52.49	50.70	43.03
Traditional MoE w/ $\mathcal{L}_{aux}$ 4E/64E	30.55	56.94	47.78	32.70	34.39	69.53	30.33	39.87	52.80	40.85	43.57
MoNE w/o $\mathcal{L}_{aux}$ 8E/64E	28.67	52.72	44.93	32.06	34.85	69.48	30.14	39.92	53.91	42.25	42.89
MoNE w/ $\mathcal{L}_{aux}$ 8E/64E	30.97	55.75	48.01	33.34	34.44	70.67	29.86	38.89	53.83	49.30	44.51

### A.3 LLM USAGE

This study utilizes Large Language Models (LLMs) to refine content, adjust formatting, construct tables, and provide writing suggestions for specific chapters.