# Autonomy-of-Experts Models

**Ang Lv** [1][2]  **Ruobing Xie** [2]  **Yining Qian** [3]  **Songhao Wu** [1]  **Xingwu Sun** [2]  **Zhanhui Kang** [2]  **Di Wang** [2]  **Rui Yan** [1]

## Abstract

Mixture-of-Experts (MoE) models mostly use a router to assign tokens to specific expert modules, activating only partial parameters and often outperforming dense models. We argue that the separation between the router's decision-making and the experts' execution is a critical yet overlooked issue, leading to suboptimal expert selection and ineffective learning. To address this, we propose Autonomy-of-Experts (AoE), a novel MoE paradigm in which experts autonomously select themselves to process inputs. AoE is based on the insight that an expert is aware of its own capacity to effectively process a token, an awareness reflected in the scale of its internal activations. In AoE, routers are removed; instead, experts pre-compute internal activations for inputs and are ranked based on their activation norms. Only the top-ranking experts proceed with the forward pass, while the others abort. The overhead of pre-computing activations is reduced through a low-rank weight factorization. This self-evaluating-then-partner-comparing approach ensures improved expert selection and effective learning. We pre-train language models having 700M up to 4B parameters, demonstrating that AoE outperforms traditional MoE models with comparable efficiency.

## 1. Introduction

Large language models (LLM) built on Mixture-of-Experts techniques (MoE, Shazeer et al., 2017; Lepikhin et al., 2021; Fedus et al., 2022) have gained increasing research and industrial attention (Jiang et al., 2024; Dai et al., 2024; Team, 2024; Sun et al., 2024). The core idea of MoE in LLMs involves dividing a large feed-forward network (FFN) into smaller FFNs, known as experts, and activating different ex-

[1]Renmin University of China [2]Machine Learning Platform Department, Tencent [3]Southeast University, China. Correspondence to: Ruobing Xie <xrbsnowing@163.com>, Rui Yan <ruiyan@ruc.edu.cn>.

*Preprint.*



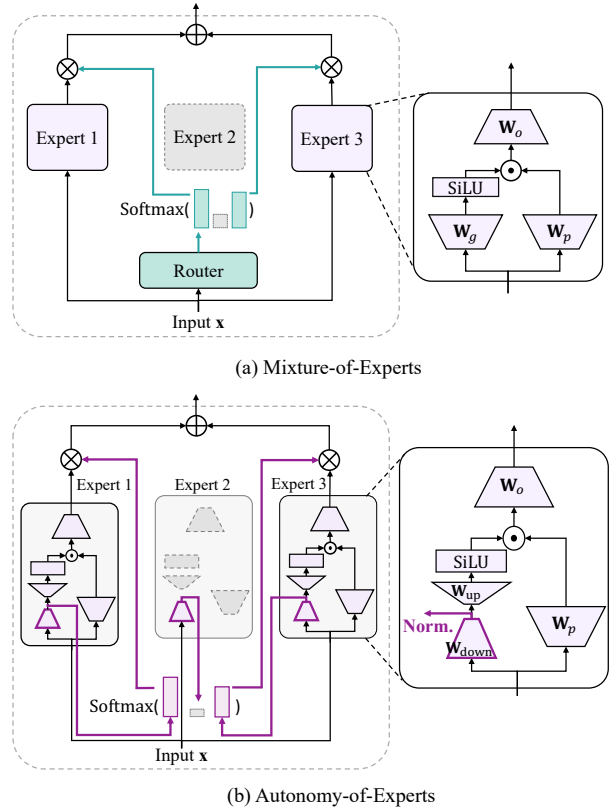(a) Mixture-of-Experts



(b) Autonomy-of-Experts

*Figure 1.* Comparison between traditional MoE and AoE. Arrows indicate data flow, while shadowed modules represent unused parameters or variables. (a) Traditional MoE models use a router to assign tokens to specific experts. This separation between the router's decision-making and the experts' execution leads to suboptimal expert selection and ineffective learning. (b) In an AoE model, experts operate autonomously. They are ranked based on their internal activation norms, and only the top-activated experts continue processing, while the others are terminated. The AoE expert architecture is modified to maintain efficiency.

perts' parameters for different inputs. The decision on which experts process which inputs are made by a router, typically an MLP-based classifier. Compared to dense models, MoE models are more efficient due to their sparse activation, and their ability to flexibly combine expert knowledge enhances downstream performance.

A critical issue in MoE is often overlooked: the separation between the router's decision-making and the experts' execution. The router cannot directly assess the experts' abilities, making its selection of an expert essentially a prediction without available labels. If the router makes an incorrect prediction, the chosen expert may struggle to process the tokens effectively, leading to increased training loss. To reduce the loss, the expert might adapt its parameters to handle these tokens, potentially conflicting with its original expertise. Alternatively, the router must learn to make better decisions through trial and error, as it lacks awareness of which experts are best suited for specific tasks, thereby wasting many training steps.

To address these challenges, we propose a novel MoE paradigm—Autonomy-of-Experts (AoE). AoE allows experts to decide whether to process inputs themselves. This design is based on the observation that experts are aware of their ability to handle inputs, an awareness reflected in the scale of their internal activations. Building on this insight, we enable all experts in an AoE layer to process every token and cache their internal activations. For each token, experts are ranked by their internal activation norms, with only the top-ranked experts continuing to process the token using the cache, while the others terminate the process. The additional overhead from caching and computations of unused experts is mitigated by factorizing the experts' weights, which compresses the inputs into low-dimensional vectors for efficient caching. Due to the autonomy of the experts, the router is eliminated. Figure 1 presents a comparative overview of traditional MoE and AoE models.

We pre-train AoE language models with up to 4 billion parameters, and they outperform traditional MoE models on downstream tasks with comparable efficiency. We provide a comprehensive analysis of AoE to highlight its advantages. These advantages include improved expert selection, more specialized experts, and more effective training, all of which contribute to better downstream performance.

## 2. Background: Mixture-of-Experts (MoE)

We focus on studies of sparse MoE models, treating each feed-forward network (FFN) module as an expert. Each FFN, or expert, is expected to possess diverse and distinct abilities, enabling the model to process inputs effectively by activating only the experts with the necessary capabilities, thereby improving efficiency. Some studies (Chen et al., 2024; Lin et al., 2024) on dense MoE do not reduce the parameter activation ratio, which is not the primary concern of this paper. In this paper, when we refer to MoE, we mean sparse MoE.

MoE-based LLMs (Jiang et al., 2024; Dai et al., 2024; Team, 2024; Lieber et al., 2024; Sun et al., 2024; Abdin et al., 2024) typically follow the FFN design in the Llama mod-

---

**Algorithm 1** A working pipeline of an MoE layer

1: **Input:** A hidden state $\mathbf{x} \in \mathbb{R}^{d_{\text{model}}}$, number of experts $n$. Initialize output $\mathbf{h} \in \mathbb{R}^{d_{\text{model}}}$ as a zero vector.
2: $\mathbf{p} = R(\mathbf{x})$          // $\mathbf{p} \in \mathbb{R}^n$
3: $\mathbf{I} = \texttt{argtopK}(\mathbf{p})$       // $\mathbf{I} \in \mathbb{R}^K$
4: $\hat{\mathbf{p}} = \texttt{Softmax}(\mathbf{p}[\mathbf{I}])$    // $\hat{\mathbf{p}} \in \mathbb{R}^K$
5: **for** $i = 1$ **to** $n$ **do**
6:     **if** $i \in \mathbf{I}$ **then**
7:        $\mathbf{h} \mathrel{+}= \hat{\mathbf{p}}[i] \cdot E_i(\mathbf{x})$
8:     **end if**
9: **end for**

---

els (Touvron et al., 2023) as an expert module. The $i$-th expert within a specific layer can be formulated as:

$$E_i(\mathbf{x}) = \left( \texttt{SiLU}(\mathbf{x}\mathbf{W}_g^i) \odot (\mathbf{x}\mathbf{W}_p^i) \right) \mathbf{W}_o^i, \qquad (1)$$

where $\mathbf{x} \in \mathbb{R}^{d_{\text{model}}}$ is the input hidden state; $\mathbf{W}_g^i, \mathbf{W}_p^i \in \mathbb{R}^{d_{\text{model}} \times d_{\text{ffn}}}$, and $\mathbf{W}_o^i \in \mathbb{R}^{d_{\text{ffn}} \times d_{\text{model}}}$ are the expert weights. This paper focuses on this classical FFN formulation.

A router (or gate) $R$ determines which expert processes which hidden state. Many studies have proposed various routing strategies, such as token choosing top experts (Shazeer et al., 2017; Lepikhin et al., 2021), expert choosing top tokens (Zhou et al., 2022; 2023), dynamic expert calls (Raposo et al., 2024; Gong et al., 2024), and refining expert selection by solving mathematical problems (Lewis et al., 2021; Clark et al., 2022), among others. Without loss of generality, our discussion focuses on token choosing the Top-$K$ experts (Shazeer et al., 2017; Lepikhin et al., 2021), but our experiments consider various strategies. Algorithm 1 presents a working pipeline of an MoE layer with a total of $n$ experts. The "$[i]$" notation in the algorithm follows Python syntax, indicating the selection of the $i$-th element in a vector or a matrix.

A challenge faced by MoE is the imbalanced expert load. MoE routers tend to disproportionately favor specific experts, resulting in suboptimal parameter utilization. Fedus et al. (2022) incorporate a load-balancing loss, controlled by a hyperparameter weight, $\alpha_{\text{aux}}$, to ensure that each expert receives a similar load for a batch $\mathcal{B}$ with $T$ tokens:

$$\mathcal{L}_{\text{aux}} = \alpha_{\text{aux}} \cdot n \cdot \sum_{i=1}^{n} \mathbf{f}_i \cdot \mathbf{P}_i, \text{where}$$

$$\mathbf{f}_i = \frac{1}{T} \sum_{\mathbf{x} \in \mathcal{B}} \mathbb{1}\left\{ i \in \texttt{argtopK}(R(\mathbf{x})) \right\}, \qquad (2)$$

$$\mathbf{P}_i = \frac{1}{T} \sum_{\mathbf{x} \in \mathcal{B}} \texttt{Softmax}(R(\mathbf{x}))[i].$$

Several variants of this auxiliary loss have been proposed (Zuo et al., 2022; Wang et al., 2024a;b; Huang et al., 2024), but they all share the common goal of maintaining a balanced load. Therefore, our discussion focuses on the balancing loss presented above.

*Table 1.* We remove routers from pre-trained MoE-LLMs and select experts during inference based on the internal activation norms of specific nodes in the computational graph. The accuracy on two challenging tasks is reported, along with the time cost (in minutes) for 8×A800-80G GPUs, which is given in parentheses. Without parameter updates, we can largely preserve accuracy under certain nodes, but this rudimentary approach requires significant improvements in efficiency.

| Node for Norm Calculation | MMLU (5-shot) | | ARC-C (5-shot) | |
|---|---|---|---|---|
| | Mixtral $8 \times$ 7B | Phi-3.5-MoE-ins. | Mixtral $8 \times$ 7B | Phi-3.5-MoE-ins. |
| $\mathbf{x}\mathbf{W}_g$ | 64.23 (42.70) | 29.43 (33.05) | 50.43 (4.40) | 28.84 (3.47) |
| $\mathbf{x}\mathbf{W}_p$ | 62.06 (42.73) | 34.60 (33.05) | 53.41 (4.40) | 40.36 (3.47) |
| $\texttt{SiLU}(\mathbf{x}\mathbf{W}_g)$ | 61.71 (43.88) | 38.03 (34.32) | 58.79 (4.51) | 47.53 (3.60) |
| $\texttt{SiLU}(\mathbf{x}\mathbf{W}_g) \odot \mathbf{x}\mathbf{W}_p$ | 66.64 (75.53) | 27.89 (52.60) | 58.79 (6.27) | 35.32 (5.42) |
| Experts' Final Outputs | 66.66 (76.15) | 29.69 (69.20) | 58.62 (7.42) | 36.35 (7.07) |
| Performance w. Router | 70.35 (24.30) | 78.20 (14.53) | 62.12 (2.50) | 67.41 (1.60) |

Several studies (Roller et al., 2021; Gururangan et al., 2022; Ren et al., 2023; Fan et al., 2020) classify tokens based on prior knowledge—such as domain, language, or hash mapping—and assign them to fixed experts. While they do not use explicit routers, they differ significantly from AoE in many respects. Most importantly, their expert selection is not determined by the experts themselves, leaving the separation between decision-making and execution unaddressed.

## 3. Method

We begin by introducing preliminary experiments that motivate the development of Autonomy-of-Experts (AoE) in Section 3.1. In Section 3.2, we refine the straightforward implementation from the preliminary experiments, improving the expert architecture to address efficiency concerns and, finally, deriving the AoE method.

### 3.1. An Insight: Experts "Know" What They Know

We present the experiment that motivated the development of AoE models.

Geva et al. (2021) interpret FFN layers as key-value memory networks, where inputs are projected into a "key" vector (e.g., $(\texttt{SiLU}(\mathbf{x}\mathbf{W}g) \odot (\mathbf{x}\mathbf{W}p))$). The "key" vector retrieves knowledge or abilities stored in the parameters through a key-value matching mechanism (e.g., multiplying by $\mathbf{W}_o$). If the experts can effectively handle the input, the "key" should be highly activated, allowing for effective retrieval. Note that this example is purely analogical; there are no defined rules to determine which internal activations behave more like the "key" and which behave more like the "value," as models are not trained with constraints that would regularize these roles.

Inspired by (Geva et al., 2021), we conducted preliminary experiments to explore whether experts in pre-trained MoE-LLMs "know" their capabilities—that is, whether the scale of their activation norms reflects their ability to handle spe-

cific inputs. Specifically, for a given pre-trained MoE-LLM, we remove all routers and let *every* expert within a layer to process each input up to a specific "pause" node in the computational graph (e.g., after $\mathbf{x}$ is multiplied by $\mathbf{W}_g$). We then ranked the experts based on the $L^2$ norm of their activations at the node.[1] The top-$K$ experts continue the forward pass from the pause node to generate the final MoE outputs, while the others are terminated. We conducted 5-shot tests on Mixtral $8 \times$ 7B (Jiang et al., 2024) and Phi-3.5-MoE-instruct (Abdin et al., 2024) using MMLU (Hendrycks et al., 2021) and ARC-Challenge (Clark et al., 2018), and investigated how much of the performance of these LLMs can be preserved using this expert-selection strategy.

Regarding which node to use for calculating the activation norm, we conducted several trials. The accuracy scores under various setups are shown in Table 1. We also report the time taken on 8×A800-80G. The test code is based on the LM Evaluation Harness (Gao et al., 2024) with a batch size of 50. Experiments across different models and tasks reveal that the optimal nodes for preserving the performance of a *pre-trained* LLM vary. This finding supports the earlier assertion that, in a *pre-trained* LLM, there is no predetermined node whose norm best reflects experts' underlying abilities. Notably, this experiment does not update any parameters and is conducted under out-of-distribution inference behavior, i.e., without routers. Despite this, performance preservation reaches up to 95% for Mixtral and 71% for Phi-3.5.

These preliminary results motivate us to train an MoE model *from scratch* with an explicit designation of the node for expert selection. We expect that the model will naturally learn to represent its awareness of its capabilities through the norm of the designated node. Such an approach could effectively address the separation between the router's decision-making and the experts' execution—a challenge inherent in traditional MoE models.

---

[1] We also evaluated the $L^1$ and $L^\infty$ norms, but these performed worse than the $L^2$ norm, as detailed in Appendix A.

**Algorithm 2** A working pipeline of an AoE layer

1: **Input:** A hidden state $\mathbf{x} \in \mathbb{R}^{d_{\text{model}}}$, number of experts $n$. Initialize output $\mathbf{h} \in \mathbb{R}^{d_{\text{model}}}$, activation cache $\mathbf{C} \in \mathbf{R}^{n \times d_{\text{low}}}$ and $\mathbf{p} \in \mathbb{R}^n$ as all zeros.
2: // In practice, we replace the following loop with a
3: // single matrix multiplication (see Eq. 4) for efficiency.
4: **for** $i = 1$ **to** $n$ **do**
5:     $\mathbf{C}[i] = \mathbf{x}\mathbf{W}_{\text{down}}^i$              // $\mathbf{C}[i] \in \mathbb{R}^{d_{\text{low}}}$
6: **end for**
7: $\mathbf{p} = \texttt{L2-Norm}(\mathbf{C}, \text{ dim=-1})$     // $\mathbf{p} \in \mathbb{R}^n$
8: $\mathbf{I} = \texttt{argtopK}(\mathbf{p})$               // $\mathbf{I} \in \mathbb{R}^K$
9: $\hat{\mathbf{p}} = \texttt{Softmax}(\mathbf{p}[\mathbf{I}])$        // $\hat{\mathbf{p}} \in \mathbb{R}^K$
10: **for** $i = 1$ **to** $n$ **do**
11:   **if** $i \in \mathbf{I}$ **then**
12:     $\mathbf{h} \mathrel{+}= \hat{\mathbf{p}}[i] \cdot \big((\texttt{SiLU}(\mathbf{C}_i\mathbf{W}_{\text{up}}^i) \odot (\mathbf{x}\mathbf{W}_p^i))\mathbf{W}_o^i\big)$
13:   **end if**
14: **end for**

## 3.2. Autonomy-of-Experts (AoE)

The following paper centers on using the norm of $\mathbf{x}\mathbf{W}_g$ to guide expert selection in our new MoE language models pre-trained from scratch. There is no technical difference or challenge in applying our method to any other node, regardless of the architecture. However, utilizing nodes other than $\mathbf{x}\mathbf{W}_g$ or $\mathbf{x}\mathbf{W}_p$ is not cost-effective.

The efficiency of the rudimentary method in Section 3.1 must be improved. The primary overhead arises from all experts computing activations for a given token, even though not all results contribute to the final MoE output. Additionally, large $d_{\text{ffn}}$-dimensional activations (14,336 for Mixtral $8 \times 7$B and 6,400 for Phi-3.5-MoE) at the pause node are cached, leading to significant memory usage.

A factorization of the $\mathbf{W}_g$ matrix can address these two issues. We decompose $\mathbf{W}_g$ into two low-rank matrices: $\mathbf{W}_{\text{down}} \in \mathbb{R}^{d_{\text{model}} \times d_{\text{low}}}$ and $\mathbf{W}_{\text{up}} \in \mathbb{R}^{d_{\text{low}} \times d_{\text{wide}}}$, where $d_{\text{low}} < d_{\text{model}} < d_{\text{wide}}$. The $i$-th AoE expert can be formulated as:

$$E_i(\mathbf{x}) = \big(\texttt{SiLU}\big(\mathbf{x}\mathbf{W}_{\text{down}}^i\mathbf{W}_{\text{up}}^i\big) \odot \big(\mathbf{x}\mathbf{W}_p^i\big)\big)\mathbf{W}_o^i, \quad (3)$$

where $\mathbf{W}_p^i \in \mathbb{R}^{d_{\text{model}} \times d_{\text{wide}}}$, and $\mathbf{W}_o^i \in \mathbb{R}^{d_{\text{wide}} \times d_{\text{model}}}$.

Algorithm 2 formulates the pipeline within an AoE layer. In each expert, $\mathbf{W}_{\text{down}}$ first compresses the input vectors into low-dimensional activations. These activations are cached as $\mathbf{C}$, and their $L^2$ norms are used to rank the experts. Given an input, the experts with the top-$K$ norms use the cache to continue the forward computation within the expert, while unchosen experts abort processing. The compressed activations significantly reduce both the cache size and the computational overhead from unselected experts. This factorization does not impair the model's expressiveness, as the weights are inherently low-rank in large language models (Li et al., 2018; Aghajanyan et al., 2021; Hu et al., 2022).

Furthermore, to enhance efficiency, the loop for calculating the activation cache (Line 2 in Algorithm 2) can be eliminated by combining the $\mathbf{W}_{\text{down}}^i$ matrices of all experts into a single large matrix. This allows the cache to be obtained through a single multiplication:

$$\begin{aligned} \hat{\mathbf{W}}_{\text{down}} &= [\mathbf{W}_{\text{down}}^1, \cdots, \mathbf{W}_{\text{down}}^n] \in \mathbb{R}^{d_{\text{model}} \times (n d_{\text{low}})} \\ \mathbf{C} &= \mathbf{x}\hat{\mathbf{W}}_{\text{down}}. \end{aligned} \quad (4)$$

The resulting $\mathbf{C} \in \mathbb{R}^{n d_{\text{low}}}$ is then reshaped into an $n \times d_{\text{low}}$ matrix for subsequent computations.

In Section 4.1, we demonstrate that an AoE model achieves up to 97% of the throughput of a traditional MoE model while also delivering superior downstream performance.

## 4. Experiments

We begin by providing a detailed analysis of our method through ablation experiments on pre-trained small language models using AoE and traditional MoE. These experiments enable us to answer key research questions related to AoE. Based on the insights gained, we scale up the language models to 4 billion parameters, demonstrating AoE's scalability.

### 4.1. Method Analysis through Small Language Models

#### 4.1.1. GENERAL SETUP

We train small language models consisting of 12 layers, each containing 12 attention heads. Each layer contains 8 experts, with the top-$K = 2$ experts selected. Models use the Llama (Touvron et al., 2023) vocabulary of size 32,000 and the same pre-RMSNorm (Zhang & Sennrich, 2019) module. We set $d_{\text{model}} = 768$ and $d_{\text{ffn}} = 3{,}072$ for traditional MoE models, while the values of $d_{\text{low}}$ and $d_{\text{wide}}$ for AoE models are variable. Specifically, in all experiments below, to ensure that the total number of parameters in an AoE model is comparable to that of an MoE model, when we adjust $d_{\text{low}}$, $d_{\text{wide}}$ is set as follows:

$$d_{\text{wide}} = \frac{3 \cdot d_{\text{model}} \cdot d_{\text{ffn}} - d_{\text{low}} \cdot d_{\text{model}}}{d_{\text{low}} + 2 \cdot d_{\text{model}}}. \quad (5)$$

The total number of model parameters is 732 million, and the number of activated parameters is 247 million.

We train models on 100 billion tokens from RedPajama (Computer, 2023), with a batch size of 4.2 million tokens, a learning rate of $2 \times 10^{-4}$, and a linear warmup over the first 4,800 steps, followed by a cosine decay schedule that reduces the learning rate to $1.28 \times 10^{-5}$ (Tow et al., 2024). The AdamW optimizer (Loshchilov & Hutter, 2019) is employed with $(\beta_1, \beta_2) = (0.9, 0.95)$, a gradient norm clipping threshold of 1, and a weight decay of 0.1.

We conduct a comprehensive evaluation of language models across a range of widely used tasks, including *ARC-easy* (Clark et al., 2018), *PIQA* (Bisk et al., 2020),

*Table 2.* Ablations were performed on 732M-parameter language models (with 247M active parameters). Each model was trained on 100 billion tokens. The results, highlighted in color, emphasize superior performance compared to configuration ②, the most common MoE setup. Bold text indicates that the configuration outperforms the best traditional MoE variant in terms of average performance.

| Configuration | ARC-E | PIQA | SIQA | WINO | HELLA | MNLI | QNLI | SST2 | AVG. |
|---|---|---|---|---|---|---|---|---|---|
| ① Traditional MoE | 39.90 | 58.43 | 35.67 | 52.09 | 27.98 | 33.09 | 49.28 | 49.66 | 43.28 |
| ②     + $\mathcal{L}_{\text{aux}}$ | 40.74 | 58.49 | 36.13 | 51.30 | 28.11 | 32.67 | 50.23 | 51.83 | 43.68 |
| ③     + $\mathcal{L}_{\text{aux}}$ + Factorized $\mathbf{W}_g$ | 40.45 | 58.65 | 36.75 | 52.09 | 28.03 | 32.55 | 50.08 | 51.03 | 43.70 |
| ④     + $\mathcal{L}_{\text{aux}}$ + Large Router | 41.41 | 57.62 | 36.64 | 52.33 | 28.34 | 33.18 | 49.53 | 50.69 | 43.71 |
| ⑤ AoE $(d_{\text{low}} = 64)$ | 39.77 | 58.71 | 35.31 | 52.33 | 28.29 | 32.78 | 50.27 | 52.98 | **43.81** |
| ⑥     + $\mathcal{L}_{\text{aux}}$ | 42.17 | 57.67 | 36.75 | 50.75 | 28.15 | 34.06 | 50.49 | 53.10 | **44.12** |
| ⑦ AoE $(d_{\text{low}} = 128)$ | 40.70 | 59.41 | 36.64 | 52.09 | 28.06 | 34.38 | 50.69 | 53.21 | **44.39** |
| ⑧     + $\mathcal{L}_{\text{aux}}$ | 41.33 | 58.65 | 36.80 | 50.75 | 28.40 | 33.71 | 49.55 | 53.10 | **44.04** |
| ⑨ AoE $(d_{\text{low}} = 256)$ | 41.08 | 58.81 | 36.44 | 51.70 | 28.23 | 32.24 | 50.54 | 53.90 | **44.12** |
| ⑩     + $\mathcal{L}_{\text{aux}}$ | 41.16 | 58.32 | 36.80 | 53.04 | 28.37 | 32.78 | 50.61 | 54.59 | **44.46** |
| ⑪ AoE $(d_{\text{low}} = 512)$ | 40.57 | 57.89 | 36.75 | 50.59 | 28.38 | 32.71 | 49.72 | 53.56 | **43.77** |
| ⑫     + $\mathcal{L}_{\text{aux}}$ | 41.16 | 57.83 | 36.75 | 52.09 | 28.30 | 34.92 | 50.67 | 50.92 | **44.08** |

*Table 3.* Comparison of traditional MoE and AoE models trained using alternative expert-selection strategies. For the Top-$P$ strategy, the number of activated parameters is input-dependent but nearly the same between the two models, whereas the expert-choice strategy activates 247 out of 732M parameters.

| Strategy | Model | ARC-E | PIQA | SIQA | WINO | HELLA | MNLI | QNLI | SST2 | AVG. |
|---|---|---|---|---|---|---|---|---|---|---|
| Top-$P$ | Traditional MoE | 41.08 | 57.96 | 37.46 | 50.36 | 28.25 | 32.79 | 50.39 | 52.64 | 43.87 |
| (Huang et al., 2024) | AoE | 41.04 | 58.65 | 36.39 | 51.07 | 28.35 | 32.96 | 51.46 | 54.36 | **44.29** |
| Expert-Choice | Traditional MoE | 40.91 | 59.09 | 37.26 | 50.75 | 28.09 | 32.11 | 50.12 | 52.75 | 43.89 |
| (Zhou et al., 2022) | AoE | 41.58 | 58.22 | 37.21 | 53.04 | 28.44 | 33.83 | 50.54 | 50.46 | **44.17** |

*SIQA* (Sap et al., 2019), *Winogrande* (Sakaguchi et al., 2019), *HellaSwag* (Zellers et al., 2019), *MNLI* (Williams et al., 2018), *MRPC* (Dolan & Brockett, 2005), *QNLI* (Wang et al., 2019), *QQP* (Wang et al., 2019), and *SST-2* (Socher et al., 2013). The first five tasks are evaluated zero-shot, while the remaining tasks are tested three-shot because models exhibit unstable performance in zero-shot scenarios, with most errors arising from incorrect answer formats. The evaluation code is based on the LM Evaluation Harness (Gao et al., 2024), and accuracy is reported in Table 2.

### 4.1.2. RESOLVING QUESTIONS REGARDING AoE

We investigate the following questions related to AoE through a series of ablation experiments.

**Question 1: How does the downstream performance of AoE compare with traditional MoE models?** We evaluated various configurations of AoE (Configs. ⑤ to ⑫) and traditional MoE models (Configs. ① to ④). Every AoE setup outperforms the best-performing MoE configuration in terms of average accuracy across eight tasks. Notably, AoE without any auxiliary loss surpasses traditional MoE models, which enhances the simplicity of training an MoE model. Additionally, AoE exhibits lower training loss, suggesting more efficient training. We elaborate on this in Question 2.
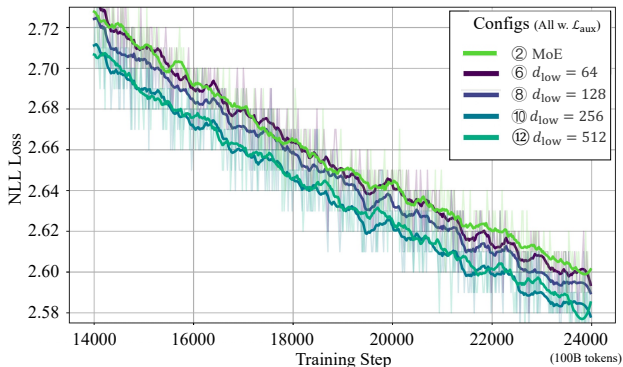


*Figure 2.* Pre-training NLL losses. All configurations shown are trained with $\mathcal{L}_{\text{aux}}$, though its value is not included in the figure.

**Question 2: What is the impact of varying $d_{\text{low}}$?** We adjusted $d_{\text{low}}$ to values of 64, 128, 256, and 512, corresponding to Configs. ⑥, ⑧, ⑩, and ⑫, respectively. The combined impact of $\mathcal{L}_{\text{aux}}$ and $d_{\text{low}}$ will be discussed in the next question. All of these variants outperform the traditional MoE model in downstream performance. The performance differences among these configurations are relatively small. The maximum performance gain occurs when $d_{\text{low}}$ is approximately one-third of $d_{\text{model}}$ (256/768). Both smaller and larger values of $d_{\text{low}}$ result in lower performance, though
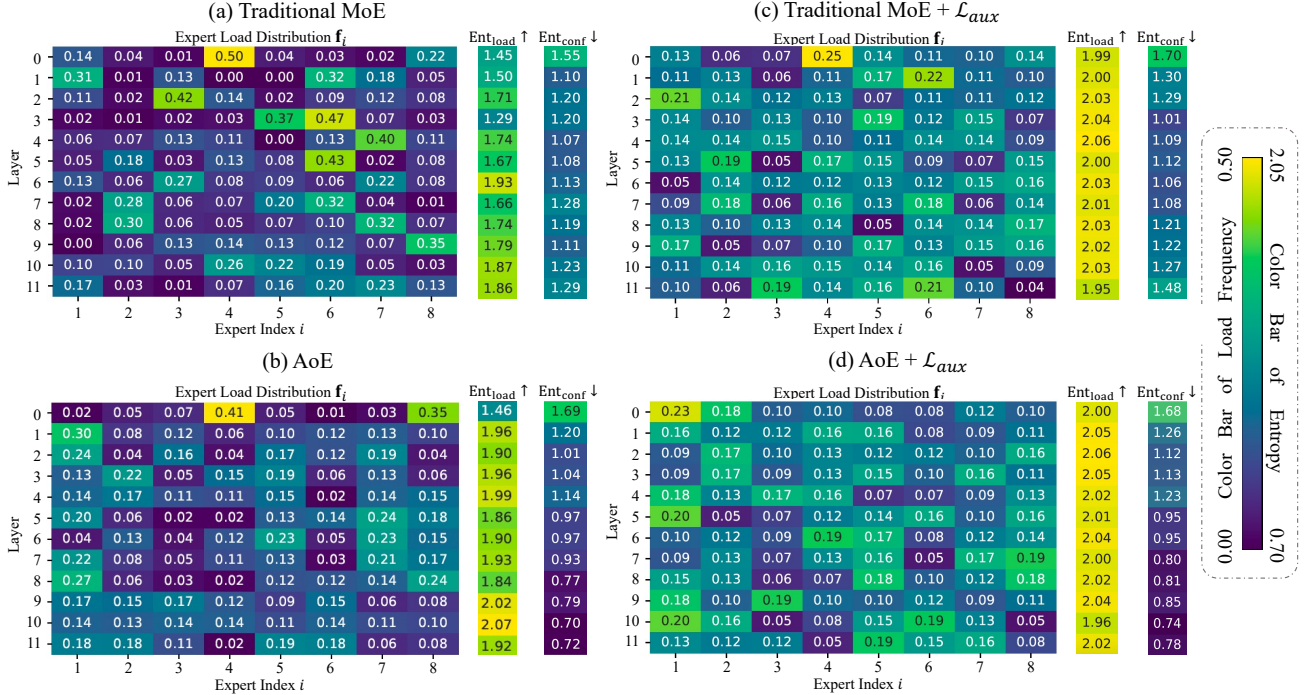
Figure 3. Statistical analysis of expert load. The figure reveals several key insights: (1) $\mathcal{L}_{\text{aux}}$ enhances load balancing in both traditional MoE and AoE. (2) AoEs generally exhibit more balanced load distributions compared to their traditional MoE counterparts, as indicated by higher $\text{Ent}_{\text{load}}$ values. (3) AoEs also demonstrate greater confidence in expert selection, reflected by lower $\text{Ent}_{\text{conf}}$ values.

they still surpass the baselines. The suboptimal performance with smaller $d_{\text{low}}$ may be due to the factorization of $\mathbf{W}_g$ into $\mathbf{W}_{\text{down}}\mathbf{W}_{\text{up}}$ being a lossy approximation when $d_{\text{low}}$ is below the true rank of $\mathbf{W}_g$. Conversely, larger $d_{\text{low}}$ introduce more noise into the activation, potentially hindering the effectiveness of the norm-based selection measure.

In Figure 2, we present the negative log-likelihood (NLL) loss during training for traditional MoE (Config. ②) and AoE models (Configs. ⑥, ⑧, ⑩, and ⑫). AoE models exhibit more effective expert learning, as evidenced by lower loss values. However, when $d_{\text{low}} = 64$ (Config. ⑥), the loss is comparable to that of traditional MoE models, suggesting that smaller $d_{\text{low}}$ values hinder AoE performance. In contrast, $d_{\text{low}} = 256$ (Config. ⑩) results in the lowest training loss overall, reinforcing the finding that setting $d_{\text{low}}$ to approximately one-third of $d_{\text{model}}$ yields the most benefits.

**Question 3: Is AoE compatible with other expert-selection strategies?** We also train language models using the Top-$P$ token-choice (Huang et al., 2024) and the Top-$K$ expert-choice strategy (Zhou et al., 2022). Table 3 shows that AoE outperforms traditional MoE models, demonstrating its generality. Details are provided in Appendix B.

**Question 4: How is the load balancing of AoE?** There are three main findings regarding load balancing.

***Finding 4.1:*** *AoE improves load balancing compared to traditional MoE models, with or without $\mathcal{L}_{\text{aux}}$.*

AoE can incorporate $\mathcal{L}_{\text{aux}}$ with minor modifications to Eq. 2, as shown below:

$$\mathcal{L}_{\text{aux}} = \alpha_{\text{aux}} \cdot n \cdot \sum_{i=1}^{n} \mathbf{f}_i \cdot \mathbf{P}_i, \text{ where}$$

$$\mathbf{f}_i = \frac{1}{T} \sum_{\mathbf{x} \in \mathcal{B}} \mathbb{1}\left\{i \in \texttt{argtopK}\left(\texttt{L2-Norm}\left(\mathbf{x}\mathbf{W}_{\text{down}}^i\right)\right)\right\}, \quad (6)$$

$$\mathbf{P}_i = \frac{1}{T} \sum_{\mathbf{x} \in \mathcal{B}} \texttt{Softmax}\left(\texttt{L2-Norm}\left(\mathbf{x}\mathbf{W}_{\text{down}}^i\right)\right)[i].$$

$\alpha_{\text{aux}}$ is determined using a validation set comprising 5 billion tokens from (Gokaslan & Cohen, 2019). Experiments indicate that $\alpha_{\text{aux}} = 0.01$ is effective for both traditional MoE and AoE models. We adopted this value across all configurations without further hyperparameter tuning.

Figure 3 illustrates expert load statistics on the SST-2 dataset (Socher et al., 2013) for Configs. ①, ② (Traditional MoE with and without $\mathcal{L}_{\text{aux}}$), ⑨, and ⑩ (AoE with and without $\mathcal{L}_{\text{aux}}$). We report both the load distribution $\mathbf{f}_i$ (as defined in Eqs. 2 and 6), representing the percentage of tokens processed by expert $i$, and the entropy of the load distribution within each layer:

$$\text{Ent}_{\text{load}} = -\sum_{i=1}^{n} \mathbf{f}_i \log \mathbf{f}_i. \quad (7)$$

Higher entropy values indicate more balanced load distributions across experts. Comparing Figures 3(a) and 3(b), without $\mathcal{L}_{\text{aux}}$, AoE achieves a more balanced load distribution in 11 out of 12 layers. Comparing Figures 3(c) and 3(d), with $\mathcal{L}_{\text{aux}}$, AoE maintains a superior overall balance. For reference, the average $\text{Ent}_{\text{load}}$ values for subfigures (c) and (d) are 2.015 and 2.023, respectively. [2]

***Finding 4.2:*** *AoE models exhibit stronger confidence in expert selection.*

We introduce the confidence entropy, denoted as $\text{Ent}_{\text{conf}}$. For each layer, we have:

$$\text{Ent}_{\text{conf}} = -\sum_{i=1}^{n} \mathbf{p}_i \log \mathbf{p}_i,$$

$$\mathbf{p}_i = \begin{cases} \texttt{Softmax}\left(\texttt{L2-Norm}\left(\mathbf{x}\mathbf{W}_{\text{down}}^i\right)\right), \text{ for AoE} \\ \texttt{Softmax}\left(R(\mathbf{x})\right), \text{ for traditional MoE} \end{cases}$$

$$(8)$$

This entropy quantifies the confidence in expert selection: lower entropy indicates a distribution closer to a one-hot vector, signifying more confident expert selection, while higher entropy reflects greater uncertainty in expert decisions. AoE exhibits significantly lower entropy, demonstrating stronger confidence in selecting experts. Furthermore, its confidence increases from shallow to deep layers, aligning with the intuitive inductive bias that shallow layers perform fundamental, non-specialized functions, whereas deeper layers handle specialized and abstract tasks (Wang et al., 2023; Lv et al., 2024). In contrast, MoE models do not display this trend, potentially suggesting more homogeneous expertise within and across layers (Wang et al., 2024a).

***Finding 4.3:*** *Beyond improved load balancing, AoE with $\mathcal{L}_{\text{aux}}$ achieves better downstream performance.*

In general, $\mathcal{L}_{\text{aux}}$ benefits both traditional MoE and AoE models. However, when $d_{\text{low}} = 128$, applying $\mathcal{L}_{\text{aux}}$ results in a decrease in accuracy, which we attribute to task-specific variations. In conclusion, as addressed in response to Question 4, AoE exhibits strong potential for advancing MoE-based LLMs, owing to its improvements in both load balancing and downstream performance.

**Question 5: Do improvements stem from the factorization of $\mathbf{W}_g$?** We examined the impact of factorizing the experts' weight matrix on performance by comparing Configurations ③ and ②. The factorization does not significantly influence performance, as expected in Section 3.2, based on findings that the weights of LLMs are inherently low-rank (Li et al., 2018; Aghajanyan et al., 2021; Hu et al., 2022). Therefore, the improvements observed with AoE are not attributed to the factorization of model weights.



*Figure 4.* Average activation norm dynamics during training. Each plot represents an expert, distinguished by color according to its layer. Experts within the same layer achieve similar activation scales, indicating that their self-evaluation criteria for determining whether they are capable of processing inputs are aligned.

**Question 6: Does the improvement of AoE come from involving more parameters in expert selection?** We increased the size of the router in MoE to include $n \cdot d_{\text{low}} \cdot d_{\text{model}}$ parameters, ensuring that the number of parameters involved in expert selection remains consistent with that of AoE models. Note that in this setup, traditional MoE models have more activated parameters in total. Comparing Config. ④ and ②, the larger router provides a slight performance benefit. However, every AoE setup still outperforms this configuration. Thus, the improvement in AoE is not primarily due to involving more parameters in expert selection.

**Question 7: How aligned are the self-evaluation criteria among experts?** In AoE models, each expert independently develops self-evaluation criteria for processing tokens, as reflected in their activation scales. This might raise concerns that some experts could become overly "egoistic," meaning their internal activations are consistently larger than those of others. For example, one expert might produce activations

---

[2] $d_{\text{low}}$ has minimal impact on the statistical metrics discussed in Question 4. As a result, we do not provide analysis for other configurations, as they offer little additional insight.
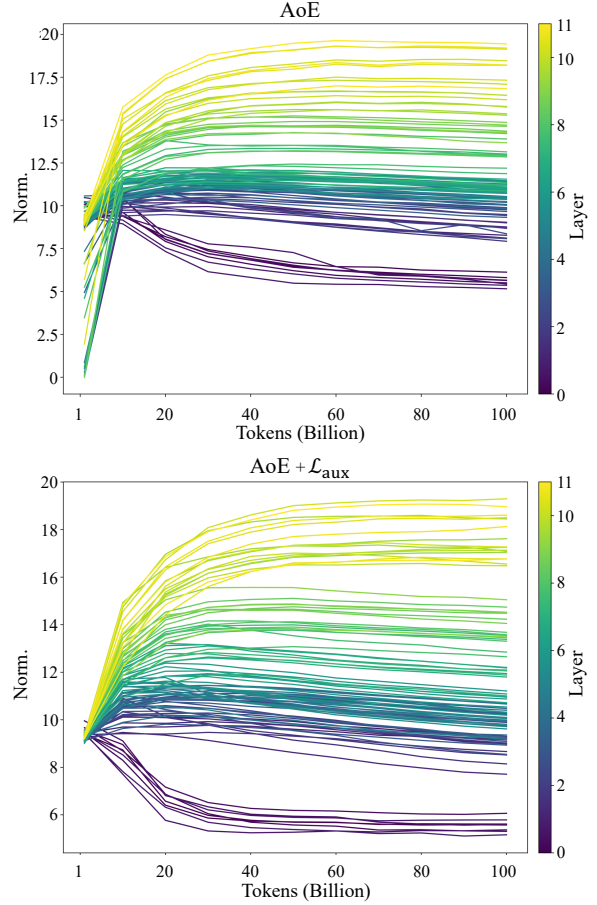
*Table 4.* For 4B-parameter LLMs (with 1.18B active parameters), AoE exhibits better downstream performance than MoE models.

| Model | ARC-E | PIQA | SIQA | WINO | HELLA | MNLI | QNLI | SST2 | AVG. |
|---|---|---|---|---|---|---|---|---|---|
| Traditional MoE | 53.70 | 65.40 | 39.10 | 51.54 | 35.80 | 32.19 | 49.77 | 57.00 | 48.06 |
| AoE | 55.98 | 65.61 | 39.87 | 52.57 | 36.77 | 35.39 | 50.05 | 61.93 | **49.80** |

with norms ranging from 10 to 20, while an "ego" expert produces activations with norms from 20 to 30, leading to biased selections that favor the "ego" expert.

We track dynamics of activation norms during pre-training. Figure 4 shows the details for Configs. ⑨ and ⑩. Except for the very initial period, experts' self-evaluation criteria are well aligned, as evidenced by clusters of same-colored plots (representing experts within the same layer). In the early stages of training without the auxiliary loss, some middle-to-upper-layer experts exhibit significantly lower activation. However, AoE naturally resolves this imbalance in activation scales during training. Alternatively, $\mathcal{L}_{\text{aux}}$ can address this imbalance earlier because it acts as a regularizer for activation norms, increasing the norm scales of underactive experts and ensuring they are used more often.

**Question 8: How Efficient is AoE?** Table 5 shows the maximum training throughput (tokens processed per second per GPU) and memory usage for both traditional MoE models and various AoE models. Here are the key findings:

***Finding 8.1:** AoE achieves up to 97% of the throughput of the traditional MoE model, with the added cost of memory.*

Additionally, note that experts in our experiments work sequentially within the same layer but in practical deployments of MoE-LLMs, experts are typically distributed across different devices and operate in parallel. Consequently, experts must wait for the most loaded expert to finish computation, resulting in idle time that can be quantified by the difference between the maximum and minimum expert loads. The total differences across layers are 1.49 for Figure 3(c) (traditional MoE) and 1.41 for Figure 3(d) (AoE). In this case, AoE can achieve an additional time reduction equivalent to processing 8% of the total tokens through a single MoE layer. Assuming an ideal load distribution where each of the 8 experts processes 12.5% of the total tokens, this reduction translates to a 64% decrease in the running time of one MoE layer. This advantage, however, is not reflected in the reported efficiency metrics.

***Finding 8.2:** In AoE, memory usage and throughput are influenced by $d_{low}$, presenting trade-offs.*

In terms of incremental memory, a smaller $d_{\text{low}}$ requires a larger $d_{\text{wide}}$, thereby increasing the memory consumption of $\mathbf{x}\mathbf{W}_{\text{up}}$ to $T \cdot d_{\text{wide}}$, where $T$ is the number of tokens. Conversely, a larger $d_{\text{low}}$ results in a larger activation cache, raising memory usage to $n \cdot T \cdot d_{\text{low}}$. For Configs. ⑥ to

*Table 5.* Throughput and memory usage comparison among several configurations. Auxiliary losses do not impact efficiency.

| Configuration | TP. (K/s) / Mem. (GB) |
|---|---|
| Traditional MoE | 51.42 / 50.61 |
| AoE $(d_{\text{low}} = 64)$ | 49.79 / 59.39 |
| AoE $(d_{\text{low}} = 128)$ | 49.42 / 57.86 |
| AoE $(d_{\text{low}} = 256)$ | 47.98 / 57.32 |
| AoE $(d_{\text{low}} = 512)$ | 46.07 / 55.90 |

⑩, $n \cdot d_{\text{low}} < d_{\text{wide}}$, making the primary memory cost stem from the larger up-projection. In contrast, Config. ⑪ and ⑫ satisfy $n \cdot d_{\text{low}} > d_{\text{wide}}$, meaning the increased memory usage is more attributable to the larger activation cache. In terms of throughput reduction, a smaller $d_{\text{low}}$ requires more computational resources for the up-projection, while a larger $d_{\text{low}}$ leads to a higher unused activation cache.

### 4.2. Pre-training Large Language Models

We pre-train LLMs with a total of 4 billion parameters, of which 1.18B are activated. The initial learning rate is $3.2 \times 10^{-4}$ (Tow et al., 2024). Each model has 24 layers, with 20 attention heads per layer. For traditional MoE models, we set $d_{\text{model}} = 1,280$ and $d_{\text{ffn}} = 5,120$. Considering the trade-offs between efficiency overhead and performance gain, we set $d_{\text{low}} = 400$ and, according to Eq. 5, derive $d_{\text{wide}} = 6,470$. Other settings follow those in Section 4.1. Both models are enhanced by $\mathcal{L}_{\text{aux}}$ with $\alpha_{\text{aux}} = 0.01$. Table 4 demonstrates that AoE outperforms traditional MoE models as they scale, with the performance improvement being more pronounced in LLMs compared to smaller models. This highlights the potential of AoE to drive advancements in larger and more powerful MoE-based LLMs.

## 5. Conclusion

We introduce Autonomy-of-Experts (AoE), a novel Mixture-of-Experts (MoE) paradigm that addresses a crucial yet widely overlooked issue: the separation between the router's decision-making and the experts' execution, which leads to suboptimal expert selection and learning. AoE selects experts based on their internal activation scales. Several architectural modifications ensure efficiency. Language models based on AoE outperform traditional MoE models in many aspects. This paper highlights the advantages of enabling MoE experts to self-select and aims to inspire the community to develop more powerful MoE-like models.

# References

Abdin, M., Aneja, J., Awadalla, H., Awadallah, A., Awan, A. A., Bach, N., Bahree, A., Bakhtiari, A., Bao, J., Behl, H., Benhaim, A., Bilenko, M., Bjorck, J., Bubeck, S., Cai, M., Cai, Q., Chaudhary, V., Chen, D., Chen, D., Chen, W., Chen, Y.-C., Chen, Y.-L., Cheng, H., Chopra, P., Dai, X., Dixon, M., Eldan, R., Fragoso, V., Gao, J., Gao, M., Gao, M., Garg, A., Giorno, A. D., Goswami, A., Gunasekar, S., Haider, E., Hao, J., Hewett, R. J., Hu, W., Huynh, J., Iter, D., Jacobs, S. A., Javaheripi, M., Jin, X., Karampatziakis, N., Kauffmann, P., Khademi, M., Kim, D., Kim, Y. J., Kurilenko, L., Lee, J. R., Lee, Y. T., Li, Y., Li, Y., Liang, C., Liden, L., Lin, X., Lin, Z., Liu, C., Liu, L., Liu, M., Liu, W., Liu, X., Luo, C., Madan, P., Mahmoudzadeh, A., Majercak, D., Mazzola, M., Mendes, C. C. T., Mitra, A., Modi, H., Nguyen, A., Norick, B., Patra, B., Perez-Becker, D., Portet, T., Pryzant, R., Qin, H., Radmilac, M., Ren, L., de Rosa, G., Rosset, C., Roy, S., Ruwase, O., Saarikivi, O., Saied, A., Salim, A., Santacroce, M., Shah, S., Shang, N., Sharma, H., Shen, Y., Shukla, S., Song, X., Tanaka, M., Tupini, A., Vaddamanu, P., Wang, C., Wang, G., Wang, L., Wang, S., Wang, X., Wang, Y., Ward, R., Wen, W., Witte, P., Wu, H., Wu, X., Wyatt, M., Xiao, B., Xu, C., Xu, J., Xu, W., Xue, J., Yadav, S., Yang, F., Yang, J., Yang, Y., Yang, Z., Yu, D., Yuan, L., Zhang, C., Zhang, C., Zhang, J., Zhang, L. L., Zhang, Y., Zhang, Y., Zhang, Y., and Zhou, X. Phi-3 technical report: A highly capable language model locally on your phone, 2024. URL https://arxiv.org/abs/2404.14219.

Aghajanyan, A., Gupta, S., and Zettlemoyer, L. Intrinsic dimensionality explains the effectiveness of language model fine-tuning. In Zong, C., Xia, F., Li, W., and Navigli, R. (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 7319–7328, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long. 568. URL https://aclanthology.org/2021. acl-long.568.

Bisk, Y., Zellers, R., Le bras, R., Gao, J., and Choi, Y. Piqa: Reasoning about physical commonsense in natural language. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7432–7439, Apr. 2020. doi: 10.1609/ aaai.v34i05.6239. URL https://ojs.aaai.org/ index.php/AAAI/article/view/6239.

Chen, Y., Lv, A., Lin, T.-E., Chen, C., Wu, Y., Huang, F., Li, Y., and Yan, R. Fortify the shortest stave in attention: Enhancing context awareness of large language models for effective tool use. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 11160–11174, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.601. URL https: //aclanthology.org/2024.acl-long.601.

Clark, A., De Las Casas, D., Guy, A., Mensch, A., Paganini, M., Hoffmann, J., Damoc, B., Hechtman, B., Cai, T., Borgeaud, S., Van Den Driessche, G. B., Rutherford, E., Hennigan, T., Johnson, M. J., Cassirer, A., Jones, C., Buchatskaya, E., Budden, D., Sifre, L., Osindero, S., Vinyals, O., Ranzato, M., Rae, J., Elsen, E., Kavukcuoglu, K., and Simonyan, K. Unified scaling laws for routed language models. In Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., and Sabato, S. (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 4057–4086. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/ v162/clark22a.html.

Clark, P., Cowhey, I., Etzioni, O., Khot, T., Sabharwal, A., Schoenick, C., and Tafjord, O. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018. URL https://arxiv.org/abs/ 1803.05457.

Computer, T. Redpajama: An open source recipe to reproduce llama training dataset, 2023. URL https://github.com/togethercomputer/ RedPajama-Data.

Dai, D., Deng, C., Zhao, C., Xu, R. X., Gao, H., Chen, D., Li, J., Zeng, W., Yu, X., Wu, Y., Xie, Z., Li, Y. K., Huang, P., Luo, F., Ruan, C., Sui, Z., and Liang, W. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models, 2024. URL https: //arxiv.org/abs/2401.06066.

Dolan, W. B. and Brockett, C. Automatically constructing a corpus of sentential paraphrases. In *Proceedings of the Third International Workshop on Paraphrasing (IWP2005)*, 2005. URL https://aclanthology. org/I05-5002.

Fan, A., Bhosale, S., Schwenk, H., Ma, Z., El-Kishky, A., Goyal, S., Baines, M., Celebi, O., Wenzek, G., Chaudhary, V., Goyal, N., Birch, T., Liptchinsky, V., Edunov, S., Grave, E., Auli, M., and Joulin, A. Beyond english-centric multilingual machine translation, 2020. URL https://arxiv.org/abs/2010.11125.

Fedus, W., Zoph, B., and Shazeer, N. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity, 2022. URL https://arxiv.org/ abs/2101.03961.

Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi, A., Foster, C., Golding, L., Hsu, J., Le Noac'h, A., Li, H., McDonell, K., Muennighoff, N., Ociepa, C., Phang, J., Reynolds, L., Schoelkopf, H., Skowron, A., Sutawika, L., Tang, E., Thite, A., Wang, B., Wang, K., and Zou, A. A framework for few-shot language model evaluation, 07 2024. URL https://zenodo.org/records/12608602.

Geva, M., Schuster, R., Berant, J., and Levy, O. Transformer feed-forward layers are key-value memories. In Moens, M.-F., Huang, X., Specia, L., and Yih, S. W.-t. (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 5484–5495, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.446. URL https://aclanthology.org/2021.emnlp-main.446.

Gokaslan, A. and Cohen, V. Openwebtext corpus. http://Skylion007.github.io/OpenWebTextCorpus, 2019.

Gong, Z., Lv, A., Guan, J., Wu, W., Zhang, H., Huang, M., Zhao, D., and Yan, R. Mixture-of-modules: Reinventing transformers as dynamic assemblies of modules. In Al-Onaizan, Y., Bansal, M., and Chen, Y.-N. (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 20924–20938, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.1164. URL https://aclanthology.org/2024.emnlp-main.1164.

Gururangan, S., Lewis, M., Holtzman, A., Smith, N. A., and Zettlemoyer, L. DEMix layers: Disentangling domains for modular language modeling. In Carpuat, M., de Marneffe, M.-C., and Meza Ruiz, I. V. (eds.), *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 5557–5576, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.407. URL https://aclanthology.org/2022.naacl-main.407.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding, 2021. URL https://arxiv.org/abs/2009.03300.

Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=nZeVKeeFYf9.

Huang, Q., An, Z., Zhuang, N., Tao, M., Zhang, C., Jin, Y., Xu, K., Xu, K., Chen, L., Huang, S., and Feng, Y. Harder tasks need more experts: Dynamic routing in moe models, 2024. URL https://arxiv.org/abs/2403.07652.

Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D. S., de las Casas, D., Hanna, E. B., Bressand, F., Lengyel, G., Bour, G., Lample, G., Lavaud, L. R., Saulnier, L., Lachaux, M.-A., Stock, P., Subramanian, S., Yang, S., Antoniak, S., Scao, T. L., Gervet, T., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mixtral of experts, 2024. URL https://arxiv.org/abs/2401.04088.

Lepikhin, D., Lee, H., Xu, Y., Chen, D., Firat, O., Huang, Y., Krikun, M., Shazeer, N., and Chen, Z. {GS}hard: Scaling giant models with conditional computation and automatic sharding. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=qrwe7XHTmYb.

Lewis, M., Bhosale, S., Dettmers, T., Goyal, N., and Zettlemoyer, L. Base layers: Simplifying training of large, sparse models. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 6265–6274. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/lewis21a.html.

Li, C., Farkhoor, H., Liu, R., and Yosinski, J. Measuring the intrinsic dimension of objective landscapes. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=ryup8-WCW.

Lieber, O., Lenz, B., Bata, H., Cohen, G., Osin, J., Dalmedigos, I., Safahi, E., Meirom, S., Belinkov, Y., Shalev-Shwartz, S., Abend, O., Alon, R., Asida, T., Bergman, A., Glozman, R., Gokhman, M., Manevich, A., Ratner, N., Rozen, N., Shwartz, E., Zusman, M., and Shoham, Y. Jamba: A hybrid transformer-mamba language model, 2024. URL https://arxiv.org/abs/2403.19887.

Lin, H., Lv, A., Chen, Y., Zhu, C., Song, Y., Zhu, H., and Yan, R. Mixture of in-context experts enhance LLMs' long context awareness. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=RcPHbofiCN.

Loshchilov, I. and Hutter, F. Decoupled weight decay regularization, 2019. URL https://arxiv.org/abs/1711.05101.

Lv, A., Chen, Y., Zhang, K., Wang, Y., Liu, L., Wen, J.-R., Xie, J., and Yan, R. Interpreting key mechanisms of factual recall in transformer-based language models, 2024. URL https://arxiv.org/abs/2403.19521.

Raposo, D., Ritter, S., Richards, B., Lillicrap, T., Humphreys, P. C., and Santoro, A. Mixture-of-depths: Dynamically allocating compute in transformer-based language models, 2024.

Ren, X., Zhou, P., Meng, X., Huang, X., Wang, Y., Wang, W., Li, P., Zhang, X., Podolskiy, A., Arshinov, G., Bout, A., Piontkovskaya, I., Wei, J., Jiang, X., Su, T., Liu, Q., and Yao, J. Pangu-sigma: Towards trillion parameter language model with sparse heterogeneous computing, 2023. URL https://arxiv.org/abs/2303.10845.

Roller, S., Sukhbaatar, S., Szlam, A., and Weston, J. E. Hash layers for large sparse models. In Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, 2021. URL https://openreview.net/forum?id=lMgDDWb1ULW.

Sakaguchi, K., Bras, R. L., Bhagavatula, C., and Choi, Y. Winogrande: An adversarial winograd schema challenge at scale, 2019. URL https://arxiv.org/abs/1907.10641.

Sap, M., Rashkin, H., Chen, D., Le Bras, R., and Choi, Y. Social IQa: Commonsense reasoning about social interactions. In Inui, K., Jiang, J., Ng, V., and Wan, X. (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 4463–4473, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1454. URL https://aclanthology.org/D19-1454.

Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., and Dean, J. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum?id=B1ckMDqlg.

Socher, R., Perelygin, A., Wu, J., Chuang, J., Manning, C. D., Ng, A., and Potts, C. Recursive deep models for semantic compositionality over a sentiment treebank. In Yarowsky, D., Baldwin, T., Korhonen, A., Livescu, K., and Bethard, S. (eds.), *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp. 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics. URL https://aclanthology.org/D13-1170.

Sun, X., Chen, Y., Huang, Y., Xie, R., Zhu, J., Zhang, K., Li, S., Yang, Z., Han, J., Shu, X., Bu, J., Chen, Z., Huang, X., Lian, F., Yang, S., Yan, J., Zeng, Y., Ren, X., Yu, C., Wu, L., Mao, Y., Xia, J., Yang, T., Zheng, S., Wu, K., Jiao, D., Xue, J., Zhang, X., Wu, D., Liu, K., Wu, D., Xu, G., Chen, S., Chen, S., Feng, X., Hong, Y., Zheng, J., Xu, C., Li, Z., Kuang, X., Hu, J., Chen, Y., Deng, Y., Li, G., Liu, A., Zhang, C., Hu, S., Zhao, Z., Wu, Z., Ding, Y., Wang, W., Liu, H., Wang, R., Fei, H., Yu, P., Zhao, Z., Cao, X., Wang, H., Xiang, F., Huang, M., Xiong, Z., Hu, B., Hou, X., Jiang, L., Ma, J., Wu, J., Deng, Y., Shen, Y., Wang, Q., Liu, W., Liu, J., Chen, M., Dong, L., Jia, W., Chen, H., Liu, F., Yuan, R., Xu, H., Yan, Z., Cao, T., Hu, Z., Feng, X., Du, D., Yu, T., Tao, Y., Zhang, F., Zhu, J., Xu, C., Li, X., Zha, C., Ouyang, W., Xia, Y., Li, X., He, Z., Chen, R., Song, J., Chen, R., Jiang, F., Zhao, C., Wang, B., Gong, H., Gan, R., Hu, W., Kang, Z., Yang, Y., Liu, Y., Wang, D., and Jiang, J. Hunyuan-large: An open-source moe model with 52 billion activated parameters by tencent, 2024. URL https://arxiv.org/abs/2411.02265.

Team, Q. Qwen1.5-moe: Matching 7b model performance with 1/3 activated parameters", February 2024. URL https://qwenlm.github.io/blog/qwen-moe/.

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G. Llama: Open and efficient foundation language models, 2023. URL https://arxiv.org/abs/2302.13971.

Tow, J., Bellagente, M., Mahan, D., and Riquelme, C. Stablelm 3b 4e1t, 2024. URL [https://huggingface.co/stabilityai/stablelm-3b-4e1t](https://huggingface.co/stabilityai/stablelm-3b-4e1t).

Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S. R. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=rJ4km2R5t7.

Wang, A., Sun, X., Xie, R., Li, S., Zhu, J., Yang, Z., Zhao, P., Han, J. N., Kang, Z., Wang, D., Okazaki, N., and zhong Xu, C. Hmoe: Heterogeneous mixture of experts for language modeling, 2024a. URL https://arxiv.org/abs/2408.10681.

Wang, K. R., Variengien, A., Conmy, A., Shlegeris, B., and Steinhardt, J. Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In *The Eleventh International Conference on Learning Representations*,

2023. URL https://openreview.net/forum?id=NpsVSN6o4ul.

Wang, L., Gao, H., Zhao, C., Sun, X., and Dai, D. Auxiliary-loss-free load balancing strategy for mixture-of-experts, 2024b. URL https://arxiv.org/abs/2408.15664.

Williams, A., Nangia, N., and Bowman, S. A broad-coverage challenge corpus for sentence understanding through inference. In Walker, M., Ji, H., and Stent, A. (eds.), *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1112–1122, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1101. URL https://aclanthology.org/N18-1101.

Zellers, R., Holtzman, A., Bisk, Y., Farhadi, A., and Choi, Y. HellaSwag: Can a machine really finish your sentence? In Korhonen, A., Traum, D., and Màrquez, L. (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1472. URL https://aclanthology.org/P19-1472/.

Zhang, B. and Sennrich, R. Root mean square layer normalization, 2019. URL https://arxiv.org/abs/1910.07467.

Zhou, Y., Lei, T., Liu, H., Du, N., Huang, Y., Zhao, V. Y., Dai, A. M., Chen, Z., Le, Q. V., and Laudon, J. Mixture-of-experts with expert choice routing. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=jdJo1HIVinI.

Zhou, Y., Du, N., Huang, Y., Peng, D., Lan, C., Huang, D., Shakeri, S., So, D., Dai, A. M., Lu, Y., Chen, Z., Le, Q. V., Cui, C., Laudon, J., and Dean, J. Brainformers: Trading simplicity for efficiency. In Krause, A., Brunskill, E., Cho, K., Engelhardt, B., Sabato, S., and Scarlett, J. (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 42531–42542. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/zhou23c.html.

Zuo, S., Liu, X., Jiao, J., Kim, Y. J., Hassan, H., Zhang, R., Gao, J., and Zhao, T. Taming sparsely activated transformer with stochastic experts. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=B72HXs80q4.

## A. Re-running Experiments in Section 3.1 Using Alternative Expert-Selection Metrics

We also use the $L^1$ and $L^\infty$ norms as expert-selection metrics in pre-trained LLMs, which resulted in poorer performance preservation compared to the $L^2$ norm. The time costs for each configuration are identical to those presented in Table 1 and are therefore omitted here for clarity. The results are shown below.

Table 6. Preliminary study results on pre-trained MoE-LLMs, selecting experts by $L^1$ norm of internal activation.

| Node for Norm Calculation | MMLU (5-shot) | | ARC-C (5-shot) | |
|---|---|---|---|---|
| | Mixtral $8 \times 7$B | Phi-3.5-MoE-ins. | Mixtral $8 \times 7$B | Phi-3.5-MoE-ins. |
| $\mathbf{x}\mathbf{W}_g$ | 51.14 | 24.15 | 41.98 | 29.01 |
| $\mathbf{x}\mathbf{W}_p$ | 39.79 | 35.87 | 40.19 | 36.35 |
| $\mathrm{SiLU}(\mathbf{x}\mathbf{W}_g)$ | 47.29 | 26.37 | 45.73 | 36.09 |
| $\mathrm{SiLU}(\mathbf{x}\mathbf{W}_g) \odot \mathbf{x}\mathbf{W}_p$ | 54.37 | 26.95 | 50.09 | 33.79 |
| Experts' Final Outputs | 57.84 | 26.56 | 52.73 | 31.31 |
| Performance w. Router | 70.35 | 78.20 | 62.12 | 67.41 |

Table 7. Preliminary study results on pre-trained MoE-LLMs, selecting experts by $L^\infty$ norm of internal activation.

| Node for Norm Calculation | MMLU (5-shot) | | ARC-C (5-shot) | |
|---|---|---|---|---|
| | Mixtral $8 \times 7$B | Phi-3.5-MoE-ins. | Mixtral $8 \times 7$B | Phi-3.5-MoE-ins. |
| $\mathbf{x}\mathbf{W}_g$ | 48.16 | 29.28 | 43.77 | 35.92 |
| $\mathbf{x}\mathbf{W}_p$ | 50.43 | 34.78 | 49.49 | 40.02 |
| $\mathrm{SiLU}(\mathbf{x}\mathbf{W}_g)$ | 54.30 | 36.38 | 47.95 | 50.85 |
| $\mathrm{SiLU}(\mathbf{x}\mathbf{W}_g) \odot \mathbf{x}\mathbf{W}_p$ | 50.72 | 26.43 | 46.08 | 33.02 |
| Experts' Final Outputs | 51.03 | 23.64 | 53.16 | 30.12 |
| Performance w. Router | 70.35 | 78.20 | 62.12 | 67.41 |

## B. Pre-training Language Models Using Alternative Expert-Selection Strategy

We train both small traditional MoE and AoE language models by assigning tokens using the Top-$P$ strategy (Huang et al., 2024). All hyperparameters and dataset details closely follow Section 4.1, except that we replace the Top-$K = 2$ strategy with Top-$P = 0.6$ following (Wang et al., 2024a). Models utilizing the Top-$P$ strategy require an additional auxiliary loss equivalent to minimizing our introduced $\mathrm{Ent}_{\mathrm{conf}}$ (Eq. 8). This ensures that the model does not learn shortcuts by assigning uniform probabilities to all experts, which would activate too many parameters to achieve lower loss. Following (Huang et al., 2024), we set the weight of this regularization term to $10^{-4}$.

Expert-choice (Zhou et al., 2022) is similar to the Top-$K$ token-choice strategy. Consider an expert-selection matrix in the shape of $T \times n$ (i.e., the router outputs in traditional MoE or the activation norms in AoE). The token-choice strategy applies the Top-$K$ operator along the $n$ dimension, whereas expert-choice applies it along the $T$ dimension. Models trained using the expert-choice strategy do not require auxiliary losses. We set the "capacity factor" to 2 (see (Zhou et al., 2022) for details), allowing each expert to process 25% of the tokens in a batch.

$d_{\mathrm{low}}$ for AoE in these two experiments is 256. Results are shown in Table 3. AoE outperforms traditional MoE models, demonstrating its generality across various expert-selection strategies.

## C. Additional Interpretation of AoE's Advantage

We provide some intuitive insights into AoE's strengths by developing a fully controlled classification task and monitoring training dynamics of both tiny AoE and MoE models. We provide details here for interested readers. This experiment is of a toy nature and not intended as a major claim or contribution.

In our setup, inputs are multivariate Gaussian vectors belonging to three classes. Classes one and two have distinct positive and negative means, respectively, while class three has a zero mean. We adjust their standard deviations to ensure no overlap within a three-sigma range. Initially, we train both tiny AoE and MoE classifiers to distinguish between classes one and two;
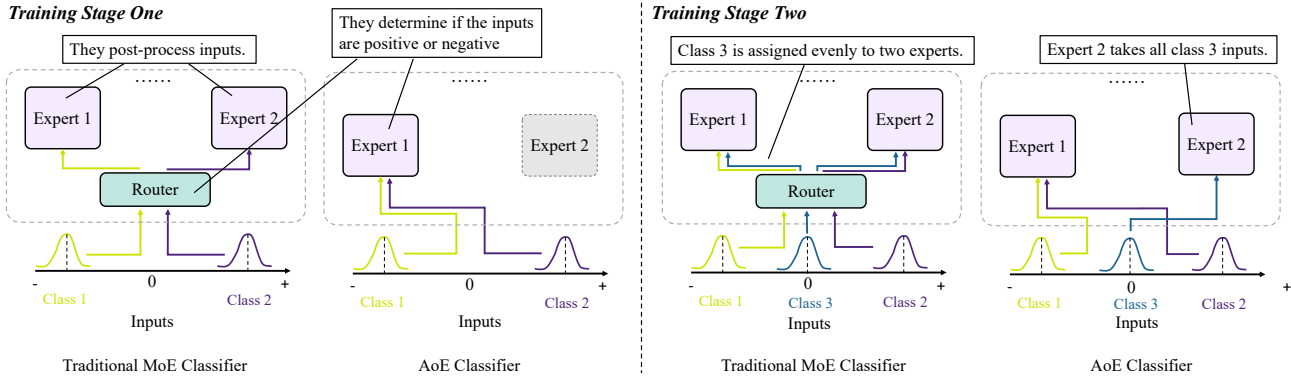
*Figure 5.* The overview of our toy experiments training tiny AoE and traditional MoE classifiers.

this is referred to as training stage one. After convergence, we introduce class three into the training process and continue training, referred to as training stage two. The classifiers consist of a single layer with two experts. Throughout training, we monitor expert behaviors, such as internal activation scales and token load. Figure 5 illustrates the pipeline and results of this toy experiment.

During training stage one, we observed that MoE classifiers assign class one and class two to different experts. This suggests that the classification role is primarily handled by the router, while the experts perform post-processing. In contrast, AoE uses only one expert to process all inputs during training stage one. Early in training, one expert identifies that the two classes are separable and develops the capability for binary classification. As training progresses, this expert's ability (reflected in increasingly larger activation norms) causes the other expert to remain naturally idle.

In training stage two, MoE evenly assigns inputs from the newly added class three to both experts. This occurs because the router has been trained for binary classification and lacks the capacity to handle out-of-distribution inputs, leading to equal prediction distribution across experts. This exacerbates the issue of homogeneous experts in the MoE classifier, as the capability to classify class three is also distributed across all experts. Conversely, in the AoE classifier, the expert handling classes one and two exhibits low activation when presented with third-class inputs. Its activation is even lower than that of the idle expert, which lacks specialization and does not resist class three inputs. As a result, the idle expert naturally handles all class three inputs. This results in heterogeneous experts within the AoE classifier: one expert manages the negative-positive classification, while the other processes zero-mean inputs.

Notably, in these toy experiments, the expert load during the first training stage is not balanced in AoE. In contrast, real-world pre-trained language models do not exhibit this imbalance, as shown in Figure 3. The reason is that the classification of input features in practical scenarios is far more complex, with a greater number of classes involved. As an evidence, when class three is added during training, AoE achieves a balanced expert load.

Comparing token assignments between the two models reveals several drawbacks of traditional MoE models:

(1) Sub-optimal expert selection: The binary classification task of distinguishing between classes one and two, which is relatively easy, could be effectively managed by a single MLP (i.e., one expert). However, MoE classifiers utilize both experts due to the router's classification behavior. This leads to under-exploitation of parameters and highlights the sub-optimal selection of experts in traditional MoE models, resulting from "the separation between the router's decision and the experts' execution."

(2) Distributed expertise: The ability to perform binary classification is distributed across two experts, preventing specialization.

The observation holds and near-zero loss is achieved as long as there is no overlap within a three-sigma range. In our experiments, we tested input dimensions and the model's $d_{\text{ffn}}$ and $d_{\text{wide}}$ parameters within the range of 32 to 256. When the input dimension is too small relative to the model dimension, the task becomes too easy to learn, and the above behavior is not observed. Conversely, if the input dimension is too large, the task becomes too difficult, preventing the loss from decreasing and rendering observed behavior uninformative.