UNION-OF-EXPERTS: EXPERTS IN MIXTURE-OF-EXPERTS ARE SECRETLY ROUTERS

Anonymous authors
Paper under double-blind review

000

001

002003004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031 032 033

034

037

038

040

041

043

044

046

047

048

050 051

052

ABSTRACT

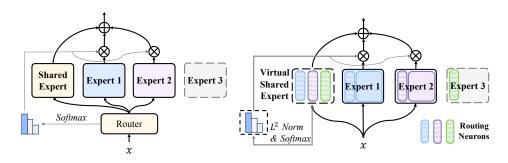
Mixture-of-Experts (MoE) is a foundational architecture in modern large language models (LLMs). However, a structural limitation has been overlooked: the router is external to the experts, rendering it unaware of their internal capabilities. This gap between routing decisions and expert capabilities limits model performance. In this paper, we demonstrate that the activations of a small subset of "routing neurons" within each routed expert's own parameters can faithfully capture the match between the expert's capabilities and input tokens. Collectively, these distributed routing neurons within each routed experts compose an implicit, capabilities-aware "router", where the norm of the routing neurons' activations suggests its corresponding expert's weight. A straightforward implementation of this design requires activating all experts to compute these routing signals, where the unselected experts' routing neurons are abandoned. To avoid the computational waste from activating unselected experts, we introduce another novel design: we unify the routing neurons of all routed experts to form a virtual shared expert, replacing the standard shared expert in MoE. In this virtual shared expert, activations are not wasted, as they serve not only for routing but also contribute to the final outputs of both the shared expert and partial of routed experts. We name this new MoE variant Union-of-Experts (UoE), drawing an analogy where the routing neuron acts as each expert's representative, and the virtual shared expert is their union, enabling the experts' autonomous selection and joint statement. We pre-train language models ranging from 1B to 3B parameters, showing that UoE consistently outperforms strong MoE baselines with comparable efficiency.

1 Introduction

Mixture-of-Experts (MoE) has garnered increasing research interest. A number of MoE-based LLMs have been proposed in recent works (DeepSeek-AI et al., 2025; Yang et al., 2025; OpenAI, 2025), exhibiting strong performance across a broad spectrum of downstream tasks. In Transformer-based MoEs, the feed-forward network (FFN) is replaced with multiple smaller expert networks, and a router dynamically routes each input token to a subset of experts. This sparse activation mechanism facilitates the training of trillion-parameter models with feasible computational overhead, establishing MoE as a fundamental architecture in modern large language models (LLMs).

However, there is a gap between routing decisions and expert capabilities. Because the router is a standalone module external to the experts, it can only infer their abilities through trial and error. When a token is inappropriately routed, the expert has to adapt to that token, compromising its specialization. To solve this, the "expert autonomy" concept has been proposed in AoE (Lv et al., 2025), wherein all experts process the token and the one with the largest activation norm (indicating the best match) is selected. While this concept improves performance, it incurs a significant computational overhead as the number of experts grows. This inefficiency contradicts the core efficiency goals of MoE models and thus limits the practical deployment of this concept in LLMs, especially under the trend of expanding total expert numbers of industrial MoE models (OpenAI, 2025; Team et al., 2025).

In this paper, we propose **Union-of-Experts** (**UoE**), a new MoE architecture that adopts the principle of expert autonomy to achieve satisfactory performance, while maintaining efficiency comparable to traditional MoE models. Figure 1 provides a comparative overview of traditional MoE and our proposed UoE architecture. The first key advancement of UoE is to *adopt only a small partial of*



(a) Mixture-of-Experts

(b) Union-of-Experts

Figure 1: A comparison of Mixture-of-Experts (MoE) and Union-of-Experts (UoE) routing mechanisms. In MoE, Expert 1 & 2 are selected based on highest router logits. In UoE, Expert 1 & 2 are selected as its routing neurons exhibit the largest activation norms. Gray modules are inactive; regions with diagonal stripe denote the routing neurons within the weight matrix.

expert neurons to represent the activation degree of an entire expert, based on a surprising finding: only a small subset of $N_s \ll D$ neurons within each expert weight, referred to **routing neurons**, is sufficient to parameterize the routing function, where D is the dimension of the intermediate activations. This reduces the computational overhead of AoE to a fraction of N_s/D . Our analysis show that the selection of routing neurons can be highly flexible. By simply pre-designating the first N_s neurons in expert's weight matrix as routing neurons before training, their activations spontaneously exhibit high correlation with those of the entire weight matrix. This indicates that these neurons can effectively represent the behavior of the majority of neurons within the expert.

Nevertheless, computing the routing neurons in each expert still introduces additional overhead. To eliminate this remaining cost, the second key advancement of UoE is to *pack routing neurons from each expert into a virtual shared expert*. This approach is grounded in a key insight: the shared expert (Dai et al., 2024) widely used in MoEs, which processes all tokens to consolidate common capabilities implicitly scattered across individual experts. UoE explicitly implements this common routing function by reusing the already-computed routing neurons, which perform the common routing function, to collectively form the output of this "virtual" shared expert. By "virtual," we mean that this is not a materialized module but a conceptual structure, describing how the outputs of routing neurons—which remain within their original experts—are reused collectively. Consequently, the computational cost of these neurons is reused rather than wasted. This allows UoE to achieve computational and memory costs identical to a standard MoE architecture while delivering superior performance.

We pre-train UoE with up to 3 billion parameters, achieving superior performance over both MoE and AoE while keeping the inference cost on par with MoE. Additionally, we present a thorough model analysis of UoE to underscore its advantages, such as improved load balance.

2 BACKGROUND AND NOTATION

2.1 MIXTURE-OF-EXPERTS (MOE)

We adopt the Gated Linear Unit (GLU) as the expert module, following mainstream MoE designs (Dai et al., 2024; Jiang et al., 2024). The *i*-th expert is parameterized by three matrices: $\mathbf{W}_g^i, \mathbf{W}_p^i \in \mathbb{R}^{d \times D}$ and $\mathbf{W}_o^i \in \mathbb{R}^{D \times d}$, with its forward pass defined as:

$$E_i(\boldsymbol{x}) = \left(\text{SiLU}(\boldsymbol{x}\boldsymbol{W}_q^i) \odot (\boldsymbol{x}\boldsymbol{W}_p^i)\right) \boldsymbol{W}_q^i. \tag{1}$$

An MoE FFN layer consists of N experts, with K experts selected to process an input token x. Adopting the design from (Dai et al., 2024), we also include a shared expert E_s that processes all tokens. This shared expert captures the common capabilities, allowing the other experts to become more specialized.

The output of an MoE FFN layer is the sum of two components: the output of a shared expert and a weighted sum of the selected expert outputs. The weights for the latter are given by a router parameterized by a matrix $\mathbf{R} \in \mathbb{R}^{d \times N}$:

$$G(\boldsymbol{x}) = \operatorname{softmax}(\boldsymbol{x}\boldsymbol{R}),$$

$$\operatorname{FFN}(\boldsymbol{x}) = \operatorname{E}_{s}(\boldsymbol{x}) + \sum_{i \in \operatorname{TopK}(G(\boldsymbol{x}))} G(\boldsymbol{x})[i] \cdot \operatorname{E}_{i}(\boldsymbol{x}).$$
(2)

2.2 AUTONOMY-OF-EXPERTS (AOE)

AoE (Lv et al., 2025) addresses the misalignment between router decisions and experts' actual capabilities by encoding the routing function $G(\boldsymbol{x})$ into the expert parameters themselves. The key insight is that the intermediate activation magnitude of an expert indicates how well its capabilities match the input token's requirements.

To reduce the computational cost associated with D-dimensional activations, \boldsymbol{W}_g^i is replaced with two low-rank matrices: $\boldsymbol{W}_{down}^i \in \mathbb{R}^{d \times r}$ and $\boldsymbol{W}_{up}^i \in \mathbb{R}^{r \times D'}$. The intermediate dimension D' is chosen to preserve the same number of parameters as the original MoE, and is given by:

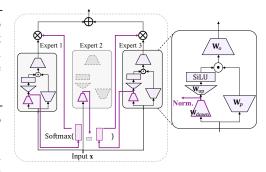


Figure 2: The overview of an AoE model.

$$D' = \frac{3Dd - dr}{r + 2d}.$$

Each token is multiplied by all \boldsymbol{W}_{down}^i matrices, and the L^2 -norms of the resulting N activations (each of dimension r) are used for expert selection. Experts with the top-K activation norms continue forward computation, while unselected experts terminate early. The routing function G and the forward pass for selected experts are defined as:

$$G(\boldsymbol{x}) = \operatorname{softmax}([g_1, g_2, \cdots, g_n]), \text{ where } g_i = \|\boldsymbol{x} \boldsymbol{W}_{down}^i\|,$$

$$E_i(\boldsymbol{x}) = \left(\operatorname{SilU}(\boldsymbol{x} \boldsymbol{W}_{down}^i \boldsymbol{W}_{up}^i) \odot (\boldsymbol{x} \boldsymbol{W}_p^i)\right) \boldsymbol{W}_o^i.$$
(3)

While AoE's autonomous expert selection leads to better downstream task performance than MoE, it introduces computational and memory overhead. The inefficiency arises because all experts compute activations, but only a fraction are used in the output. This waste scales with an increased N and a decreased K. Therefore, this paper focuses on achieving autonomous selection with an efficiency comparable to vanilla MoE, independent of N and K.

3 METHODOLOGY

3.1 MOTIVATION

To improve efficiency, AoE introduces factorization of W_g . Paradoxically, this design traps AoE in a dilemma: it must contend with either substantial computational overhead or excessive memory access. Consequently, factorization itself becomes the fundamental bottleneck to further efficiency-wise advancement in AoE. Our following analysis reveals this inherent dilemma. The detailed derivation of the results in this subsection can be found in Appendix A.

We show that AoE introduces additional FLOPs per token, which grow linearly with the factorization rank r compared to a vanilla MoE (with identical parameter count) as:

$$\Delta \text{FLOPs} = 2 \cdot d \cdot r \cdot (N - K). \tag{4}$$

Additionally, AoE incurs extra memory overhead (per token) given by:

$$\Delta \text{Mem} = \max(Nr, 4K(D'-D)). \tag{5}$$

We visualize AoE's computational and memory overhead as a function of r in Figure 3. The results clearly show that regardless of the value of r, AoE is bounded by either memory or computational resources. A rank r between 48 and 80 offers a relatively more favorable trade-off: although it still incurs significant memory overhead, the computational cost is substantially reduced. However, for wide models with large d and D, setting r this low leads to unstable training, rendering this theoretically optimal range impractical.

This dilemma motivates a new realization of autonomous expert selection, which for practicality and scalability must improve efficiency by eliminating the root cause of waste rather than relying on low-rank factorization.

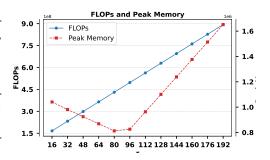


Figure 3: The efficiency dilemma introduced by factorization: for any rank r, AoE is inevitably bounded by either computation or memory.

3.2 ROUTING NEURONS ACCELERATE AUTONOMOUS ROUTING

Model structure Through extensive trials, we identified a promising approach that successfully maintains autonomous expert selection based on activation norms while achieving high efficiency without relying on factorization. As no factorization is applied, each expert in our model, namely UoE, is parameterized identically to a vanilla MoE (Eq. 1) using standard dense weight matrices.

We find that only a small subset of neurons within each expert's weight matrix is sufficient to parameterize the routing function. We refer to these as *routing neurons*. Notably, the selection of these routing neurons proves highly flexible (refer to Appendix B for more details). UoE operates by simply pre-designating the first $N_s \ll D$ neurons of each expert weight matrix as routing neurons before training. These neurons, being part of an expert's parameters, are marked with a tilde superscript:

$$\widetilde{\boldsymbol{W}_{q}^{i}} = \boldsymbol{W}_{q}^{i}[:,:N_{s}], \quad \widetilde{\boldsymbol{W}_{p}^{i}} = \boldsymbol{W}_{p}^{i}[:,:N_{s}], \quad \widetilde{\boldsymbol{W}_{o}^{i}} = \boldsymbol{W}_{o}^{i}[:N_{s},:],$$

For any input x, UoE performs autonomous expert selection based on the activation intensity (measured by L^2 norm) of routing neurons. This approach is motivated by prior work (Lv et al., 2025; Geva et al., 2021) which establishes that high activation magnitude indicates a module is well-aligned with the input. Another fundamental premise of UoE is that the activation magnitude of the routing neurons is highly correlated with that of their entire expert, a correlation we show in Section 4.4 is spontaneously reinforced during training.

Formally, we define the routing function G in UoE as:

$$G(\boldsymbol{x}) = \operatorname{softmax} \left(\operatorname{TopK} \left[g_1, g_2, \cdots, g_n \right] \right), \text{ where}$$

$$g_i = \| \operatorname{SiLU}(\boldsymbol{x} \widetilde{\boldsymbol{W}_q^i}) \odot (\boldsymbol{x} \widetilde{\boldsymbol{W}_p^i}) \|.$$

$$(6)$$

Because these routing neurons separately located in each routed expert collaboratively function as an "autonomous routing function", UoE, like AoE, eliminates the separate, explicit router module.

3.3 VIRTUAL SHARED EXPERT IMPROVES ACTIVATION UTILIZATION EFFICIENCY

We observe that routing neurons, activated on every token, functionally resemble a shared expert, which processes all tokens regardless of which experts are selected or not. We therefore consolidate them into a virtual shared expert, which replaces the conventional shared expert. This ensures the contributions of routing neurons from unselected experts are not wasted, fundamentally resolving the inherent computation and memory inefficiencies of AoE models. By "virtual," we mean that during training, these neurons are not physically restructured into a single module but remain within their original experts; their consolidation is an abstract concept describing how their activations are collectively reused beyond mere routing.

To be specific, the virtual shared expert consists of three virtual matrices during training:

$$W_g^s = \left(\begin{array}{c|c} \widetilde{\boldsymbol{W}_g^1} & \widetilde{\boldsymbol{W}_g^2} & \dots & \widetilde{\boldsymbol{W}_g^N} \end{array} \right),$$

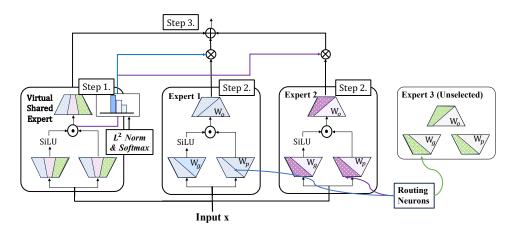


Figure 4: In UoE, the first N_s neurons in each parameter matrix are designated as routing neurons. These neurons process every token, with their activations used to compute routing logits. During training, while these neurons remain distributed across experts, they collectively function as a virtual shared expert—their outputs contribute to the final prediction like a standard shared expert, regardless of whether their host expert is selected. During inference, this virtual expert is materialized as a single module. The forward pass of UoE consists of three steps: (1) computing the activations of the routing neurons to obtain routing logits (also obtaining the output of the virtual shared expert), (2) performing expert routing using the routing logits and activating selected routed experts, and (3) merging the outputs of the virtual shared expert with those of the routed experts.

$$\begin{split} W_p^s &= \left(\begin{array}{c|c} \widetilde{\boldsymbol{W}_p^1} & \widetilde{\boldsymbol{W}_p^2} & \dots & \widetilde{\boldsymbol{W}_p^N} \end{array} \right), \\ W_o^s &= \left(\begin{array}{c|c} \widetilde{\boldsymbol{W}_o^1}^\top & \widetilde{\boldsymbol{W}_o^2}^\top & \dots & \widetilde{\boldsymbol{W}_o^N}^\top \end{array} \right)^\top. \end{split}$$

We set the number of routing neurons per parameter matrix as $N_s = \text{round}(D/K)$, matching the parameter count of a standard shared expert. This ensures that UoE has identical memory and computational overhead to a conventional MoE with the same N and K.

During inference, the virtual shared expert is materialized as a single module, ensuring UoE's checkpoint compatibility with all well-developed kernels designed for accelerating standard MoE models.

4 EXPERIMENTS

4.1 MAIN RESULTS AND ANALYSIS

General Setup. We pre-train language models with 1B parameters to verify the effectiveness of UoE. Our language model consists of 8 Transformer layers. For each Transformer layer, we employ the multi-head attention mechanism with a total of 8 attention heads. We substitute all FNN layers with MoE layers while keeping the number of expert activations consistent across all methods. The MoE baseline is configured with a shared expert following the setup in (Dai et al., 2024).

We pre-train our language models with 100B tokens from FineWeb datasets (Penedo et al., 2024), and use the Llama tokenizer for tokenization. For training setups, we employ the AdamW optimizer with $(\beta_1, \beta_2) = (0.9.0.95)$, a gradient norm clipping threshold of 1, and weight decay as 0.1. We use a learning rate of 1×10^{-3} with 1000 steps linear warmup, followed by a cosine decay scheduler.

We evaluate these language models across 8 widely used benchmarks, including ARC (Clark et al., 2018), PIQA (Bisk et al., 2020), HellaSwag (Zellers et al., 2019), SCIQ (Welbl et al., 2017), Winogrande (Sakaguchi et al., 2019), MNLI (Wang et al., 2018), QNLI (Wang et al., 2018) and RTE (Wang et al., 2018). These benchmarks assess the models' capabilities in language understanding, question answering, and natural language inference. All evaluations are performed using the LM Evaluation

Table 1: Results for the validation experiments on 1B parameter language models. We compare models with different numbers of activated experts, both with and without the auxiliary load balancing loss. Colored entries highlight improvements over the MoE baseline, while bold text mark the best results within each experimental setting.

Model	Num.	\mathcal{L}_{aux}	ARC-E	PIQA	HELLA	SCIQ	WINO	MNLI	QNLI	RTE	AVG.
MoE	8	✓	62.54	68.88	36.74	81.60	52.49	32.78	51.04	49.46	54.44
AoE	8	/	64.60	69.59	36.62	83.30	51.22	34.13	50.01	48.86	54.79
UoE	8	✓	63.09	69.64	37.07	82.40	52.88	33.89	50.05	51.50	55.07
MoE	8	X	62.75	68.23	36.62	81.10	51.85	33.12	49.95	50.18	54.23
AoE	8	X	62.29	68.17	36.32	82.20	54.14	33.71	49.78	49.10	54.46
UoE	8	X	64.56	69.10	36.86	81.50	52.09	33.02	49.91	49.46	54.56
MoE	4	✓	61.45	67.52	35.27	77.10	50.75	33.25	49.83	46.45	52.70
AoE	4	/	61.57	68.61	36.07	82.40	52.01	33.12	49.80	50.30	54.24
UoE	4	1	62.25	68.66	35.67	81.70	54.70	33.62	50.20	48.98	54.47

Harness (Gao et al., 2024). The first five tasks are evaluated zero-shot. For the remaining three tasks, we report their average performance under 0-shot, 3-shot and 5-shot to reduce randomness.

Experimental Results. We present the main results in Table 1. We pre-train 1B-parameter language models with varying number of expert activation, both with and without the auxiliary load-balancing loss. UoE consistently outperforms both MoE and AoE models in overall performance across all of these configurations, which further demonstrates the effectiveness of UoE's model design (note that UoE is more efficient than AoE, refer to Section 4.2 for more detailed discussions).

Notably, UoE achieves more substantial performance improvements under a sparser expert activation setting (activating 3 out of 64 experts), which is a defining characteristic of modern

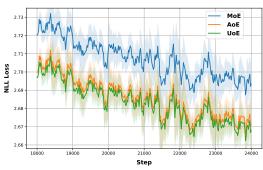


Figure 5: Pre-training NLL loss comparison.

MoE architectures. It implies that UoE could better select effective expert combinations among larger numbers of routed experts. Figure 5 illustrates the pre-training negative negative log-likelihood (NLL) loss of UoE and baseline methods in this setup. UoE exhibits a lower training loss during the pre-training phase, indicating its higher efficiency in parameter updates.

4.2 EFFICIENCY ANALYSIS OF UOE

In this section, we analyze the efficiency of UoE in comparison with the baseline methods, focusing primarily on (1) training efficiency metrics and (2) expert loading balance.

TFLOPs, Peak Memory and Throughput. We begin by conducting a comparative analysis of UoE's training efficiency. Table 2 reports the training achieved TFLOPs, peak memory usage and throughput of UoE and baseline methods during pre-training.

We observe that UoE achieves a 19.8% improvement in training throughput over AoE while maintaining downstream performance that is better or competitive with AoE, and superior to MoE. Meanwhile, UoE incurs computational

Table 2: Achieved training TFLOPs, Memory and Throughput.

	TFLOPs	Mem.(GB)	TP. (K/s)
MoE	90.40	63.93	604.00
AoE	78.29	71.51	509.00
UoE	86.51	63.96	610.00

overhead that is nearly identical to MoE at inference time. Consequently, we contend that UoE is an

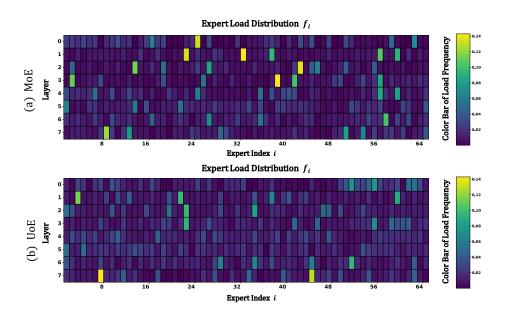


Figure 6: Expert Loading Distribution of UoE and MoE.

efficient implementation for expert autonomy and resolve the dilemma of AoE without compromising on effectiveness.

Load Balance of UoE. The imbalanced expert load is a critical challenge leading to the computational overhead of MoEs (Fedus et al., 2022). Prior study shows that AoE achieves better load balancing than traditional MoE.

We compare UoE with MoE to investigate whether UoE can enhance load balancing in the absence of an auxiliary loss. Specifically, we sample 1,000 instances from Wikitext-2 (Merity et al., 2016) as a calibration set and examine their expert loading patterns. Figure 6 visualizes the expert loading for our pre-trained UoE alongside MoE, where the load distribution f_i for the i-th expert on a batch of T tokens is defined as:

$$\mathbf{f}_{i} = \frac{1}{T} \sum_{\mathbf{x} \in \mathcal{B}} \mathbb{1} \left\{ i \in \operatorname{argtopK} \left(G \left(\mathbf{x} \right) \right) \right\}.$$

Except for the final layer, UoE achieves consistently better load balance, with far fewer cases of the imbalance observed in the shallow layers of MoE. Table 3 compares layer-wise entropy of expert selection of MoE and UoE to highlight their differences. The results shows that UoE maintains a more balanced expert load across nearly all layers, even without an auxiliary load-balancing loss.

Table 3: Entropy of expert selection. Higher entropy indicates more balanced expert loads.

Ent _{load} ¹	Ent _{load}	Ent ³ _{load}	Ent _{load} ⁴	Ent ⁵ _{load}	Ent _{load}	Ent _{load} ⁷	Ent _{load} ⁸
3.45 3.70	3.23 3.62						3.42 3.31

4.3 ABLATION STUDIES

We take our pre-trained UoE with the auxiliary balancing loss, keeping 8 experts activated as the basic setup, to evaluate the effectiveness of UoE's various designs.

Ablation Study of the Virtual Shared Expert. We perform ablation experiments to valid its contribution. We first highlight that the virtual shared expert is crucial in pre-trained UoE models. For configurations ① and ②, we deactivate different experts in the pre-trained language model and

Table 4: Analysis of model designs in UoE through ablation studies.

	Configuration	ARC-E	PIQA	HELLA	SCIQ	WINO	MNLI	QNLI	RTE	AVG.
	UoE	63.09	69.64	37.07	82.40	52.88	33.89	50.05	51.50	55.07
1 2	w.o shared expert-v1w.o shared expert-v2	53.83 62.42	66.21 69.48	33.80 37.16	75.80 81.90	50.36 52.09	33.93 33.76	50.27 49.86	51.50 51.74	51.96 54.80
3	w.o shared expert-v3	65.19	69.53	36.67	81.60	49.88	33.52	50.05	49.22	54.46
4	double N_s	63.72	68.28	36.58	84.20	51.30	34.11	49.97	50.66	54.85
⑤ ⑥ ⑦	$\mathbf{x}\mathbf{W}_p$ $\mathbf{x}\mathbf{W}_g$ SiLU $(\mathbf{x}\mathbf{W}_g)$	63.72 63.97 63.51	70.08 69.21 69.48	36.69 37.25 36.76	82.50 80.70 82.40	51.70 52.09 53.35	32.95 33.62 33.73	50.00 51.06 49.58	50.18 49.58 49.22	54.73 54.69 54.75

observe its downstream performance changes. In configuration ①, we disable the virtual shared expert and activate only the routed experts; In configuration ②, we always keep the virtual shared expert active and reduce the number of activated experts to ensure a fair comparison. Given this, we find that the shared expert exerts a significant impact on downstream performance. Configuration ① demonstrates substantially inferior performance compared with ②. This indicates that the shared expert in UoE truly learns abilities compulsory that the routed experts have not captured.

We also pre-train UoE without the virtual shared expert from scratch. In configuration ③, the routing neurons are not reactivated and are used simply for expert routing. As shown in Table 4, the absence of the activated shared expert once again leads to a decline in model performance.

Ablations Study on the Selection of N_s . We perform ablation studies to investigate the effect of varying N_s , the number of routing neurons. Specifically, we double the number of routing neurons and pre-train the model from scratch. This setting will increases an extra shared experts, while the number of activated routed experts is reduced to keep the total count of active experts constant. We do not explore alternative settings, as they would result in an excessive number of shared experts. Our results show that even doubling the number of routing neurons does not improve model performance and may even cause a slight degradation in capability.

Ablation Study of Expert Selection Strategies. By default, we use the activation intensity of the i-th expert for routing decisions, which is formulated in Equation 6. For configurations 5 to 7, we attempt to use the L^2 norm of other intermediate nodes within the computation graph for routing.

We pre-train these variants from scratch and present their downstream performance in Table 4, together with the nodes used for norm calculation. The results show that these variants achieve overall performance comparable to the default configuration ($\mathtt{SiLU}(\mathbf{x}\mathbf{W}_g)\odot\mathbf{x}\mathbf{W}_p$), albeit slightly lower.. Their training time is nearly identical. Overall, these results justify the use of activation intensity.

4.4 Consistency of Expert Selection between Expert and Routing Neuron Activations

We argue that the activation patterns of routing neurons closely reflect those of their corresponding experts. To clarify that, we perform expert routing based on the experts' activation intensity, rather than the routing neurons' in a pre-trained UoE model. To be specific, we activate the top-k experts with the highest activation values, and directly evaluate UoE's downstream performance without further training. Table 5 presents the results, and only a minor performance drop is observed.

Table 5: Performance change when using experts' activation intensity instead of routing neurons.

Model	Act.	ARC-E	PIQA	HELLA	SCIQ	WINO	MNLI	QNLI	RTE AVG.
UoE	Neurons	63.09	69.64	37.07	82.40	52.88	33.89	50.05	51.50 55.07
	Expert	61.49	68.72	36.49	82.10	51.54	34.25	50.16	50.30 54.38

Table 6: For 3B-paramter LLMs, UoE exhibits consistent downstream performance. Colored entries show improvements over the MoE baseline; bold text indicates the best results.

Model	ARC-E	ARC-C	PIQA	HELLA	SCIQ	WINO	AVG.
MoE	63.64	31.48	70.62	39.52	89.40	51.22	57.65
AoE UoE	64.44 69.07	31.57 33.11	70.24 73.18	40.34 41.96	88.80 87.10	53.35 52.80	58.12 59.54

4.5 VALIDATION OF UOE WITH LARGER MODEL SIZE

We pre-train UoE and its competitors with a total of 3 billion parameters. We follow most of the architectural settings from Section 4.1. For these 3B-parameter language models, each model consists of 20 layers and 20 attention heads, with the hidden dimension expanded to 1280. The number of experts is kept consistent with the previous setup, and 7 routed experts are activated. We adjust training parameters accordingly to better suit the training.

At larger parameter scales, UoE consistently outperforms MoE and AoE models, with improvements becoming increasingly pronounced as the model size grows. This highlights the potential of scaling UoE to even greater parameter sizes to further boost its capabilities.

5 RELATED WORK

Mixture-of-Experts. The Mixture-of-Experts (MoE) paradigm was originally proposed as a modular neural network framework in which a gating function assigns inputs to specialized experts (Jacobs et al., 1991; Jordan & Jacobs, 1994). More recently, MoE has been integrated into large-scale Transformers to achieve trillion-parameter models with sparse computation (Lepikhin et al.; Fedus et al., 2022). Subsequent work has focused on improving efficiency through balanced expert assignment (Lewis et al., 2021) and system-level optimizations for distributed training (Hwang et al., 2023; Gale et al., 2022). Despite these advances, sparse MoE models continue to face challenges such as routing instability and expert redundancy. To mitigate these issues, DeepSeekMoE (Dai et al., 2024) introduces shared experts, which provide stable coverage of common knowledge while routed experts focus on specialization. In addition, its fine-grained expert partitioning further enhances efficiency and encourages more diverse expert behaviors. In this work, we adopt most of the configurations from DeepseekMoE. In contrast, our virtual shared expert is constructed from all routing neurons, thereby functioning both as the shared expert and as the mechanism for autonomous routing.

Expert Selection Strategies. Prior work on expert selection has explored a variety of routing mechanisms to determine which experts to activate from a set of N candidates. Top-k routing (Lepikhin et al.) activates a fixed number of experts per token based on router-assigned scores, while Top-p routing dynamically selects experts until a cumulative probability threshold p is reached. Despite these differences, most approaches rely on a centralized router to assign tokens to experts. In contrast, Lv et al. (2025) eliminates the router entirely by allowing experts to self-activate, thereby achieving expert selection in a fully decentralized manner. In this paper, we improve AoE's expert autonomy by addressing efficiency issues and replacing low-rank factorization with routing neurons.

6 Conclusion

In this paper, we introduce UoE, a novel MoE variant that perform expert autonomy routing. UoE leverages only a small subset of neurons in each expert to capture the expert's overall activation, effectively addressing the efficiency challenges encountered in previous work. Moreover, we treat these routing neurons collectively as a shared expert to further enhance activation utilization efficiency. We hope that UoE can inspire the community to pursue more effective autonomy-based routing strategies to mitigate the decoupling between routing decisions and expert capabilities.

ETHICS STATEMENT

486

487 488

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

This work focuses on the development of a Mixture-of-Experts (MoE) model. Our study does not involve human subjects, personally identifiable information, or sensitive data. We do not foresee any direct ethical or societal risks arising from our methodology or experiments.

REPRODUCIBILITY STATEMENT

We have added our code to the supplementary materials, and all the data used is open-source. The experimental setup is detailed in Section 4.1. Unless noted, all experiments use the same settings. Overall, these practices make our results reproducible.

REFERENCES

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

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*, 2018.

Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli Luo, Chong Ruan, Zhifang Sui, and Wenfeng Liang. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models, 2024. URL https://arxiv.org/abs/2401.06066.

DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2025. URL https://arxiv.org/abs/2412.19437.

William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: scaling to trillion parameter models with simple and efficient sparsity. J. Mach. Learn. Res., 23(1), January 2022. ISSN 1532-4435.

- Trevor Gale, Deepak Narayanan, Cliff Young, and Matei Zaharia. Megablocks: Efficient sparse training with mixture-of-experts, 2022. URL https://arxiv.org/abs/2211.15841.
 - Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model evaluation harness, 07 2024. URL https://zenodo.org/records/12608602.
 - Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 5484–5495, 2021.
 - Changho Hwang, Wei Cui, Yifan Xiong, Ziyue Yang, Ze Liu, Han Hu, Zilong Wang, Rafael Salas, Jithin Jose, Prabhat Ram, Joe Chau, Peng Cheng, Fan Yang, Mao Yang, and Yongqiang Xiong. Tutel: Adaptive mixture-of-experts at scale, 2023. URL https://arxiv.org/abs/2206.03382.
 - Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. Adaptive mixtures of local experts. *Neural Computation*, 3:79–87, 1991. URL https://api.semanticscholar.org/CorpusID:572361.
 - Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024. URL https://arxiv.org/abs/2401.04088.
 - Michael I. Jordan and Robert A. Jacobs. Hierarchical mixtures of experts and the em algorithm. *Neural Comput.*, 6(2):181–214, March 1994. ISSN 0899-7667. doi: 10.1162/neco.1994.6.2.181. URL https://doi.org/10.1162/neco.1994.6.2.181.
 - Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. In *International Conference on Learning Representations*.
 - Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. Base layers: Simplifying training of large, sparse models. In *International Conference on Machine Learning*, pp. 6265–6274. PMLR, 2021.
 - Ang Lv, Ruobing Xie, Yining Qian, Songhao Wu, Xingwu Sun, Zhanhui Kang, Di Wang, and Rui Yan. Autonomy-of-experts models. *arXiv preprint arXiv:2501.13074*, 2025.
 - Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models, 2016.
 - OpenAI. gpt-oss-120b & gpt-oss-20b model card, 2025. URL https://arxiv.org/abs/2508.10925.
 - Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin A Raffel, Leandro Von Werra, Thomas Wolf, et al. The fineweb datasets: Decanting the web for the finest text data at scale. *Advances in Neural Information Processing Systems*, 37:30811–30849, 2024.
 - Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *arXiv preprint arXiv:1907.10641*, 2019.
 - Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong, Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei Gao, Hongcheng Gao, Peizhong Gao, Tong Gao, Xinran Gu, Longyu Guan, Haiqing Guo, Jianhang Guo, Hao Hu, Xiaoru Hao, Tianhong He, Weiran He, Wenyang He, Chao Hong, Yangyang Hu, Zhenxing Hu, Weixiao Huang, Zhiqi Huang, Zihao Huang, Tao Jiang, Zhejun Jiang, Xinyi Jin,

 Yongsheng Kang, Guokun Lai, Cheng Li, Fang Li, Haoyang Li, Ming Li, Wentao Li, Yanhao Li, Yiwei Li, Zhaowei Li, Zheming Li, Hongzhan Lin, Xiaohan Lin, Zongyu Lin, Chengyin Liu, Chenyu Liu, Hongzhang Liu, Jingyuan Liu, Junqi Liu, Liang Liu, Shaowei Liu, T. Y. Liu, Tianwei Liu, Weizhou Liu, Yangyang Liu, Yibo Liu, Yiping Liu, Yue Liu, Zhengying Liu, Enzhe Lu, Lijun Lu, Shengling Ma, Xinyu Ma, Yingwei Ma, Shaoguang Mao, Jie Mei, Xin Men, Yibo Miao, Siyuan Pan, Yebo Peng, Ruoyu Qin, Bowen Qu, Zeyu Shang, Lidong Shi, Shengyuan Shi, Feifan Song, Jianlin Su, Zhengyuan Su, Xinjie Sun, Flood Sung, Heyi Tang, Jiawen Tao, Qifeng Teng, Chensi Wang, Dinglu Wang, Feng Wang, Haiming Wang, Jianzhou Wang, Jiaxing Wang, Jinhong Wang, Shengjie Wang, Shuyi Wang, Yao Wang, Yejie Wang, Yiqin Wang, Yuxin Wang, Yuzhi Wang, Zhaoji Wang, Zhengtao Wang, Zhexu Wang, Chu Wei, Qianqian Wei, Wenhao Wu, Xingzhe Wu, Yuxin Wu, Chenjun Xiao, Xiaotong Xie, Weimin Xiong, Boyu Xu, Jing Xu, Jinjing Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu, Xinran Xu, Yangchuan Xu, Ziyao Xu, Junjie Yan, Yuzi Yan, Xiaofei Yang, Ying Yang, Zhen Yang, Zhilin Yang, Zonghan Yang, Haotian Yao, Xingcheng Yao, Wenjie Ye, Zhuorui Ye, Bohong Yin, Longhui Yu, Enming Yuan, Hongbang Yuan, Mengjie Yuan, Haobing Zhan, Dehao Zhang, Hao Zhang, Wanlu Zhang, Xiaobin Zhang, Yangkun Zhang, Yizhi Zhang, Yongting Zhang, Yu Zhang, Yutao Zhang, Yutong Zhang, Zheng Zhang, Haotian Zhao, Yikai Zhao, Huabin Zheng, Shaojie Zheng, Jianren Zhou, Xinyu Zhou, Zaida Zhou, Zhen Zhu, Weiyu Zhuang, and Xinxing Zu. Kimi k2: Open agentic intelligence, 2025. URL https://arxiv.org/abs/2507.20534.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 353–355, Brussels, Belgium, November 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-5446. URL https://aclanthology.org/W18-5446.

Johannes Welbl, Nelson F. Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. In Leon Derczynski, Wei Xu, Alan Ritter, and Tim Baldwin (eds.), *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pp. 94–106, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/W17-4413. URL https://aclanthology.org/W17-4413/.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine really finish your sentence? In Anna Korhonen, David Traum, and Lluís Màrquez (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/.

A DISCUSSION OF THE TRADE-OFF BETWEEN MEMORY CONSUMPTION AND COMPUTATIONAL OVERHEAD

For simplification, we omit the cost of the router and the FLOPs of a traditional MoE layer is:

$$FLOPs = 3 \cdot TK (2D \cdot d).$$

For an arbitrary AoE layer, the FLOPS it requires is:

$$FLOPs = 2 \cdot TK (2D' \cdot d) + TK (2D' \cdot r) + TN (2d \cdot r),$$

where D' is the FFN hidden size of AoE to ensure the same number of parameters as MoE as:

$$D' = \frac{3 \cdot D \cdot d - d \cdot r}{r + 2 \cdot d}.$$

Compared with MoE, AoE introduces an overhead of FLOPs that is:

$$\Delta FLOPs = 2 T \cdot d \cdot r \cdot (N - K). \tag{7}$$

B TRIALS ON SELECTING ROUTING NEURONS WITHIN EXPERTS

Motivated by our preliminary explorations, we investigate the idea of fixing a subset of neurons as routing neurons to enable expert autonomy. In FFNs, neurons are dynamically activated based on input. Despite that, our goal is to identify a subset of key neurons that effectively capture the overall activation pattern. Our initial approach dynamically selects important neurons during training and then fixes this subset during inference, allowing dominant weights in the experts' parameters to be located on the fly. A simple strategy uses the L^2 -norm to identify high-contributing neurons. Although this incurs higher training cost than standard MoE, it remains substantially more efficient than AoE. To further improve efficiency, we explore whether permanently fixing neurons could work. Our experiments further confirm its feasibility.

C THE USE OF LARGE LANGUAGE MODELS (LLMS)

This paper employed an LLM solely to refine our manually written draft, including improving word choice, grammar correctness, and sentence fluency.