

GROVE MOE: TOWARDS EFFICIENT AND SUPERIOR MOE LLMs WITH ADJUGATE EXPERTS

Anonymous authors

Paper under double-blind review

ABSTRACT

The Mixture of Experts (MoE) architecture is a cornerstone of modern state-of-the-art (SOTA) large language models (LLMs). MoE models facilitate scalability by enabling sparse parameter activation. However, traditional MoE architecture uses homogeneous experts of a uniform size, activating a fixed number of parameters irrespective of input complexity and thus limiting computational efficiency. To overcome this limitation, we introduce Grove MoE, a novel architecture incorporating experts of varying sizes, inspired by the heterogeneous big.LITTLE CPU architecture. This architecture features novel adjugate experts with a dynamic activation mechanism, enabling model capacity expansion while maintaining manageable computational overhead. Building on this architecture, we present GroveMoE-Base and GroveMoE-Inst, 33B-parameter LLMs developed by applying an upcycling strategy to the Qwen3-30B-A3B-Base model during mid-training and post-training. GroveMoE models dynamically activate 3.14-3.28B parameters based on token complexity and achieve performance comparable to SOTA open-source models of similar or even larger size.

1 INTRODUCTION

Recent advancements in Large Language Models (LLMs) have spurred the adoption of the Mixture of Experts (MoE) architecture (Jiang et al., 2024; Yang et al., 2025; Liu et al., 2024; Google DeepMind, 2025; Huo et al., 2025) considering its great model capacity. The MoE model operates by dynamically routing input tokens to a relevant subset of multiple experts, enhancing computational efficiency and scalability.

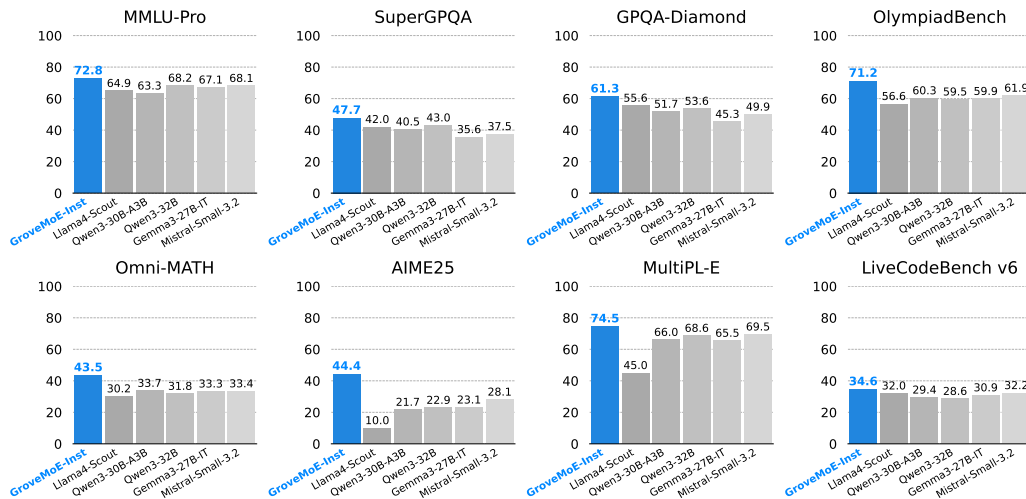


Figure 1: The performance of GroveMoE-Inst[‡] and its counterparts. GroveMoE-Inst achieves performance comparable to open-source SOTA LLMs of similar or even larger scales.

054 However, a key limitation of conventional MoE models is their reliance on homogeneous experts,
055 which activates a fixed number of parameters irrespective of input token complexity. Given that to-
056 ken complexity varies, computational resources should ideally be allocated dynamically with more
057 resources for complex tokens and fewer for simple ones (Huang et al., 2024). The current rigid-
058 ity precludes such fine-grained control over computation. Inspiration from the big.LITTLE CPU
059 architecture (Greenhalgh, 2011), which uses heterogeneous cores to manage computational load ef-
060 ficiently, we employ experts of varying sizes within MoE models could enable dynamic resource
061 allocation based on computational demand.

062 Drawing inspiration from the structure of trees, we introduce Grove MoE, a novel architecture that
063 leverages parallel adjudgate experts to expand model capacity efficiently. Grove reflects our archi-
064 tectural design, where experts are organized into disjoint groups. Similar to how branches share a
065 common trunk within a tree cluster, experts within a Grove group share a single adjudgate expert.
066 If multiple activated experts belong to the same group, their shared adjudgate expert is computed
067 only once before being added to each expert’s output. This mechanism enables a form of dynamic
068 computation allocation, allowing Grove MoE to expand model capacity with high computational
069 efficiency. Crucially, the Grove MoE architecture is compatible with existing routing mechanisms,
070 eliminating the need for complex, manually designed routing strategies (Huang et al., 2024; Wang
071 et al., 2024b) to manage expert activation counts. Moreover, we apply a loss-free expert loading
072 balance strategy during mid-training (Liu et al., 2024; Su, 2025).

073 Furthermore, the upcycling strategy (Komatsuzaki et al., 2023; Nakamura et al., 2025) has been
074 widely used for upcycling dense models, which offer larger model capacity. Recent studies (Baykal
075 et al., 2023; Wu et al., 2024; Chen et al., 2025) have also shown that expanding model capabilities
076 through parallel computation, especially by reusing existing weights, is an effective strategy com-
077 parable to direct parameter scaling. Consequently, we develop our GroveMoE architecture based on
078 pre-trained MoE models through mid-training (Meta-AI, 2025; Yang et al., 2025) and post-training
079 stages. We summarize our main contributions as follows:

- 080 • We introduce the Grove MoE architecture, which features a new mechanism of dynamic
081 computation allocation, allowing for parameter expansion while maintaining manageable
082 computational costs.
- 083 • We develop GroveMoE-Base and GroveMoE-Inst, open-source models that dynamically
084 activate 3.14-3.28B parameters, upcycled based on Qwen3-30B-A3B-Base through mid-
085 training and post-training.
- 086 • We conduct empirical evaluations showing that our GroveMoE models achieve excellent
087 performance comparable to open-source SOTA LLMs of similar or even larger scales.

088 2 RELATED WORKS

089
090 **big.LITTLE Architecture.** The big.LITTLE CPU architecture (Greenhalgh, 2011) offers a com-
091 pelling model for computational efficiency by integrating high-performance big cores and energy-
092 efficient LITTLE cores within a single processor, dynamically routing tasks to the appropriate core
093 type. In contrast, traditional MoE architectures (Yang et al., 2025; Liu et al., 2024) typically employ
094 homogeneous experts of uniform size, analogous to a processor with only one type of core, which
095 leads to suboptimal efficiency. Drawing inspiration from the big.LITTLE architecture, we propose
096 an MoE architecture where experts vary in computational capacity and the experts are dynamically
097 activated. In this paper, we introduce the Grove MoE architecture, which materializes this concept
098 through novel adjudgate experts and a dynamic activation mechanism.

099 **MoE Architecture with Dynamic Activation.** Prior research has explored dynamic activation of
100 expert counts in MoE models to mitigate the ineffectiveness of fixed top- k routing for modeling
101 targets of different complexity. Naive approaches (Jin et al., 2024; Zeng et al., 2024) indirectly vary
102 the active expert count by including blank or constant experts in the routing pool. DynMoE (Guo
103 et al., 2024) enables a top-any gating mechanism to choose any number of experts. ReMoE (Wang
104 et al., 2024b) activates experts with positive scores via the ReLU function. Both DynMoE and Re-
105 MoE face the challenge of needing explicit mechanisms to manage the upper bound on the number
106 of activated experts so as to avoid potential high computation overhead. Moreover, their relatively
107 complex routing strategies are incompatible with the well-established top- k routing mechanisms,
which brings potential issues in practice. In contrast, our GroveMoE layer intrinsically achieves

dynamic activation by assigning adjugate experts to separate groups. This approach guarantees a controllable activation count and thus manageable computation overhead. It requires no specialized router modifications, and the excellent compatibility makes it widely applicable. We provide more comparisons in Appendix A.1.

Upcycling Strategy. The performance of LLMs is intrinsically related to the model capacity. Various upcycling methods (Komatsuzaki et al., 2023; Nakamura et al., 2025) have been proposed to upcycle dense models, leveraging knowledge from pre-training models. This paper develops GroveMoE-Base and GroveMoE-Inst models by applying the upcycle strategy to the MoE model Qwen3-30B-A3B-Base (Yang et al., 2025).

3 ARCHITECTURE

3.1 TRADITIONAL MOE LAYER

An MoE layer comprises n experts $\{E_i\}_{i=1}^n$ and a router R . Let $\mathbf{x} \in \mathbb{R}^d$ represent the feature of an input token, where d denotes the feature dimension. The routing scores are calculated as $\boldsymbol{\rho} = R(\mathbf{x}) \in \mathbb{R}^n$ and the output of the i -th expert is $\mathbf{e}_i = E_i(\mathbf{x}) \in \mathbb{R}^d$. The final output of the MoE layer for each token is a weighted sum of the outputs from a selected subset of experts:

$$\mathbf{y} = \sum_{i \in \arg \text{top}k(\boldsymbol{\rho})} \rho_i \mathbf{e}_i, \quad (1)$$

where $\arg \text{top}k(\cdot)$ selects the indices of the top k routing scores.

3.2 GROVE MOE WITH ADJUGATE EXPERTS

Expanding model capacity during mid-training can preserve existing knowledge while providing additional resources to acquire complex skills. However, under the upcycling strategy, directly duplicating each expert’s parameters would disturb the original distribution of $\boldsymbol{\rho}$. Specifically, when experts are duplicated, k must be scaled proportionally by the same scale factor. However, k should be controlled to avoid introducing significant activation parameters. Alternatively, extending the parameters of each expert E_i maintains the original $\boldsymbol{\rho}$ distribution and offers a viable solution.

The AltUp architecture (Baykal et al., 2023) demonstrates that introducing parallel blocks can increase model capacity without incurring substantial computational overhead. Moreover, scaling parallel computation by reusing existing parameters has proven effective for enhancing model capability (Wu et al., 2024; Chen et al., 2025). Inspired, we introduce Grove MoE, a novel architecture that leverages parallel adjugate experts for efficient model capacity expansion. In the Grove MoE layer, we divide n experts, $\{E_i\}_{i=1}^n$, into g disjoint groups, where g divides n , and each group contains n/g experts. The groups $\{G_j\}_{j=1}^g$ can be defined as follows:

$$G_j = \left\{ E_i \mid \left\lfloor \frac{i-1}{n/g} \right\rfloor + 1 = j \right\}, \text{ for } j = 1, 2, \dots, g, \quad (2)$$

where $\lfloor \cdot \rfloor$ denotes the floor function. Motivated by big.LITTLE architecture (Greenhalgh, 2011), we introduce an adjugate expert A_j for each group G_j . Notably, the capacity of the adjugate expert A_j can differ from that of the expert E_i . The modified output $\bar{\mathbf{e}}_i$ is then computed as:

$$\bar{\mathbf{e}}_i = E_i(\mathbf{x}) + \lambda A_j(\mathbf{x}), \text{ where } j = \left\lfloor \frac{i-1}{n/g} \right\rfloor + 1. \quad (3)$$

Here, λ is the scaling factor for the adjugate expert. The final output of the Grove MoE layer is:

$$\mathbf{y} = \sum_{i \in \arg \text{top}k(\boldsymbol{\rho})} \rho_i \bar{\mathbf{e}}_i = \sum_{i \in \arg \text{top}k(\boldsymbol{\rho})} \rho_i (E_i(\mathbf{x}) + \lambda A_j(\mathbf{x})), \text{ where } j = \left\lfloor \frac{i-1}{n/g} \right\rfloor + 1. \quad (4)$$

The key advantage of the Grove MoE architecture is dynamic computation allocation. To illustrate, consider experts E_r and E_s where $\left\lfloor \frac{r-1}{n/g} \right\rfloor = \left\lfloor \frac{s-1}{n/g} \right\rfloor$. According to Equation (4):

$$\begin{aligned} \rho_r \bar{\mathbf{e}}_r + \rho_s \bar{\mathbf{e}}_s &= \rho_r (E_r(\mathbf{x}) + \lambda A_j(\mathbf{x})) + \rho_s (E_s(\mathbf{x}) + \lambda A_j(\mathbf{x})) \\ &= \rho_r E_r(\mathbf{x}) + \rho_s E_s(\mathbf{x}) + (\rho_r + \rho_s) \lambda A_j(\mathbf{x}), \text{ where } j = \left\lfloor \frac{r-1}{n/g} \right\rfloor + 1. \end{aligned} \quad (5)$$

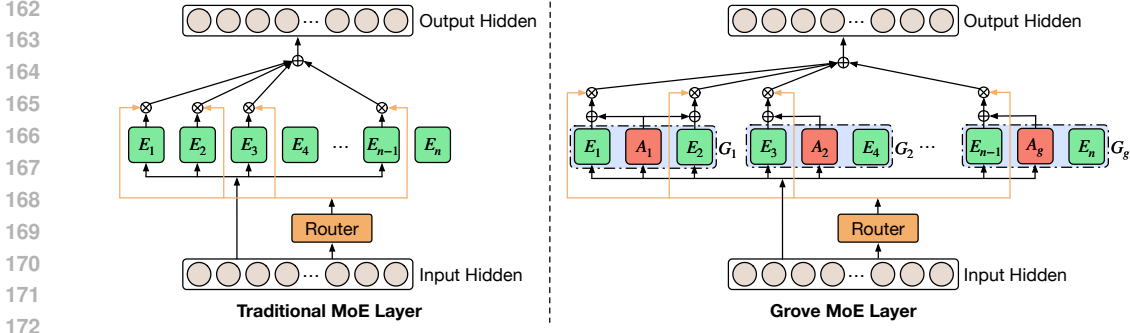


Figure 2: Comparison between the traditional MoE layer and our Grove MoE layer with dynamic computation allocation. For clarity, we configure $n/g = 2$ and $k = 4$ for the Grove MoE layer.

For the adjugate expert A_j , the routing weight $\rho_A = (\rho_r + \rho_s)\lambda$ should be restricted to be no more than the weight of experts E_r and E_s , especially for the upcycling mid-training case. Generally, we restrict $\lambda \leq 1.0/(n/g) = g/n$. For instance, for a MoE model with $n=128$ experts and $g=64$ groups, we restrict $\lambda \leq 0.5$.

As shown in Figure 2, if multiple activated experts belong to the same group, their adjugate expert A_j is computed only once before being added to each expert’s output, scaled by the sum of their routing weights. **Meanwhile, if all experts are not activated, their adjugate expert is also not computed.** According to Equation (5), the number of activated adjugate expert A_j ranges from $\lceil \frac{k}{n/g} \rceil$ to k for each Grove MoE layer, where $\lceil \cdot \rceil$ denotes the ceiling function. This mechanism enables a form of dynamic computation allocation, allowing Grove MoE to expand model capacity with high computational efficiency.

3.3 EXPERTS LOADING BALANCE

In MoE models, workload imbalance among experts can induce routing collapse and reduce computational efficiency. To better balance load distribution while maintaining model performance, we employ an auxiliary-loss-free load balancing strategy (Liu et al., 2024). Specifically, we introduce a dynamic expert-wise bias term \mathbf{b} to adjust routing scores ρ . The gating mechanism in Equation (4) is therefore modified as:

$$\mathbf{y} = \sum_{i \in \arg \text{top}k(\rho + \mathbf{b})} \rho_i \bar{\mathbf{e}}_i, \quad (6)$$

where \mathbf{b} is updated via sign gradient descent to minimize imbalance.

Formally, we define $\mathbf{F} = \mathbb{E}(\mathbf{f})$ as the current load distribution across experts under the bias \mathbf{b} , where $\mathbf{f} = [f_1, f_2, \dots, f_n]$ and f_i denotes the assignment probability for each token:

$$f_i = \begin{cases} 1/k, & i \in \arg \text{top}k(\rho + \mathbf{b}), \\ 0, & \text{otherwise.} \end{cases} \quad (7)$$

The uniform load distribution is $\mathbf{Q} = [\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n}]$. To optimize load balancing (Su, 2025), \mathbf{b} is updated as:

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \frac{\mathbf{F} - \mathbf{Q}}{\sqrt{\frac{1}{n} \sum_{i=1}^n (F_i - Q_i)^2}}. \quad (8)$$

Specifically, Equation (8) uses RMSNorm to normalize the imbalance $(\mathbf{F} - \mathbf{Q})$ for better workload balance. To address the sensitivity of the parameter α in Equation (8), resulting from its coupling with the Sigmoid router (Liu et al., 2024; Su, 2025), we decouple the routing calculation. The final output \mathbf{y} is computed as:

$$\mathbf{y} = \sum_{i \in \arg \text{top}k(\rho^{(\sigma)} + \mathbf{b})} \rho_i^{(h)} \mathbf{e}_i, \quad (9)$$

Table 1: Architecture exploration on different expert group settings. The highest and second-best scores are shown in bold and underlined, respectively.

Expert Groups	-	-	$g = 64; h = 128$	$g = 32; h = 256$	$g = 16; h = 512$
Architecture	MoE	MoE	Grove MoE	Grove MoE	Grove MoE
# Total Params	30B	30B	33B	33B	33B
# Activated Params	3B	3B	3.14B-3.28B	3.14B-3.57B	3.14B-4.11B
# Avg. Activation	3B	3B	3.26B	3.51B	3.93B
# Training Tokens	0B	50B	50B	50B	50B
MMLU	81.58	<u>82.56</u>	82.57	82.12	80.23
CMMLU	80.63	86.54	<u>86.47</u>	85.54	85.57
SuperGPQA	36.10	35.93	36.22	36.09	<u>36.12</u>
GSM8K	89.39	89.54	90.07	90.83	<u>90.67</u>
MATH	59.75	65.90	66.52	<u>66.60</u>	66.64
GPQA-Diamond	39.39	38.38	<u>41.41</u>	39.39	44.95
HumanEval+	83.54	83.54	85.37	<u>84.75</u>	84.15
MBPP+	71.96	74.34	75.66	<u>75.13</u>	74.07
Average	67.79	69.59	70.54	70.06	<u>70.30</u>

Table 2: Architecture exploration on different scaling factors. The highest and second-best scores are shown in bold and underlined, respectively.

Scaling Factor	-	-	$\lambda = 0.20$	$\lambda = 0.10$	$\lambda = 0.05$
Architecture	MoE	MoE	Grove MoE	Grove MoE	Grove MoE
# Total Params	30B	30B	33B	33B	33B
# Training Tokens	0B	50B	50B	50B	50B
MMLU	81.58	82.56	79.69	<u>82.57</u>	82.62
CMMLU	80.63	<u>86.54</u>	86.33	86.47	86.55
SuperGPQA	36.10	35.93	33.99	<u>36.22</u>	36.32
GSM8K	89.39	89.54	<u>90.14</u>	90.07	90.83
MATH	59.75	65.90	66.62	<u>66.52</u>	65.86
GPQA-Diamond	39.39	38.38	<u>43.43</u>	41.41	43.94
HumanEval+	83.54	83.54	85.37	85.37	84.76
MBPP+	71.96	74.34	75.93	<u>75.66</u>	75.13
Average	67.79	69.59	70.19	<u>70.54</u>	70.75

where $\rho^{(h)}$ is the output of a Softmax router and $\rho^{(\sigma)}$ is the output of a Sigmoid router:

$$\rho_i^{(h)} = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \text{ and } \rho_i^{(\sigma)} = \frac{1}{1 + e^{-x_i}}. \quad (10)$$

This decoupled approach enables using a constant value of $\alpha = 0.001$, as used in Liu et al. (2024).

3.4 REUSE OF PRE-TRAINED WEIGHTS

Building on the concept of upcycling strategy (Komatsuzaki et al., 2023; Nakamura et al., 2025; Wu et al., 2024), we leverage pre-trained weights from MoE models. During initialization of our Grove MoE architecture, each expert E_i is derived from a pre-trained MoE layer. To ensure structural coherence, other components, such as the normalization and attention layers, are directly copied from the pre-trained transformer block. Additionally, the down-projection blocks of newly inserted modules $\{A_j\}_{j=1}^g$ are zero-initialized. The remaining weights in $\{A_j\}_{j=1}^g$ follow a normal distribution with a standard deviation of 0.006 (Liu et al., 2024).

4 MID-TRAINING

4.1 MID-TRAINING SETTING

The mid-training stage is designed to target specific proficiencies, such as reasoning and code generation. The corpus utilized for this stage is a diverse collection of textual and non-textual data, encompassing sources including web content, books, academic papers, mathematics, programming code, etc. In total, the high-quality corpus consists of approximately 400 billion tokens.

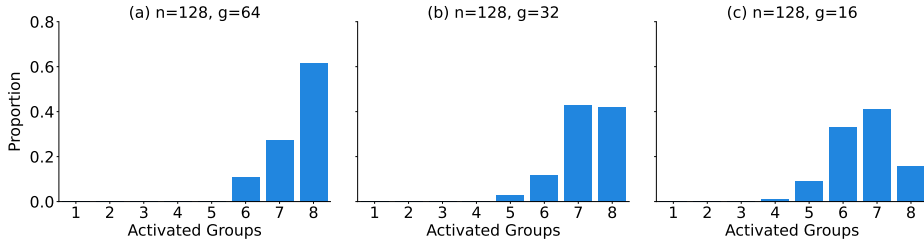


Figure 3: The group routing distribution across three configurations with different group numbers g . As the number of groups decreases, the average number of adjugate experts activations decreases.

Moreover, our comprehensive evaluation of the base models developed via mid-training assesses five core capabilities: general knowledge, scientific knowledge, reasoning, mathematics, and coding. The details of the evaluation benchmarks for mid-training are illustrated in Appendix B.1.

4.2 ARCHITECTURE EXPLORATION

We explored several key design choices for the model architecture, focusing on the configuration of the expert groups and the scaling factor in our Grove MoE architecture. The architecture of the shared parallel experts $\{A_j\}_{j=1}^g$ adopts the design established in Qwen3 MoE architecture (Yang et al., 2025). Architectural exploration is conducted using 50B tokens sampled from the mid-training dataset, with evaluations performed across diverse task types. The exploration is based on the Qwen3-30B-A3B-Base model (Yang et al., 2025), using direct mid-training without upcycling as the baseline.

Expert Groups. As detailed in Table 1, we evaluate the impact of expert group configurations on model performance while maintaining approximately 33B total parameters. Three configurations are compared: (1) $g = 64, h = 128$; (2) $g = 32, h = 256$; and (3) $g = 16, h = 512$, where h denotes the intermediate dimension of $\{A_j\}_{j=1}^g$. For general knowledge, language understanding, and code generation, configuration with $g = 64, h = 128$ achieves optimal performance. Meanwhile, a configuration with $g = 16, h = 512$ yields superior results for complex mathematical reasoning and STEM tasks. Notably, three configurations outperform the baseline in terms of average performance, demonstrating the effectiveness of our Grove MoE architecture for expanding model capacity.

Scaling Factor. Table 2 evaluates the influence of the expert output scaling factor λ . For general knowledge and language understanding, $\lambda = 0.05$ achieves peak performance. For mathematics and STEM tasks, both $\lambda = 0.05$ and $\lambda = 0.20$ excel, while $\lambda = 0.20$ is optimal for code generation. In general, Grove MoE with a smaller scaling factor outperforms the baseline across multiple tasks.

Group Routing Analysis. To analyze the group routing distribution, we sample 1 million tokens from the mid-training dataset. For an LLM with $n = 128$ experts, Figure 3 illustrates the distribution across three configurations. Notably, when analyzing the routing patterns across diverse task domains, including STEM, coding, and general knowledge, the distributions remain largely consistent and similar to the overall pattern illustrated in Figure 3. With a large number of groups ($g = 64$), expert activation is broadly distributed, with most experts assigned to 7–8 groups. In contrast, configurations with fewer groups ($g = 32$ and $g = 16$) exhibit highly consolidated expert activation. This consolidation directly impacts computational efficiency. With the $g = 64$ configuration, the average number of activated parameters is 3.26B, reducing computation by approximately 5%. As the number of groups decreases, the computational savings increase. For the $g = 16$ configuration, the savings reach approximately 20%.

4.3 HYPER-PARAMETERS

The GroveMoE-Base model is trained based on Qwen3-30B-A3B-Base (Yang et al., 2025). We employ the AdamW optimizer (Loshchilov & Hutter, 2017) with $\beta_1 = 0.9, \beta_2 = 0.95$, weight decay of 0.1, and gradient clipping at 1.0. Training uses a maximum sequence length of 8192 tokens, with the model trained on 400 billion tokens at a batch size of 16 million tokens. During the mid-training stage, we employ a cosine learning rate scheduler to decay the learning rate from

Table 3: Comparison among GroveMoE-Base and other strong open-source baselines. The highest and second-best scores are shown in bold and underlined, respectively.

	Mistral-Small-3.1 Base-2503	Gemma3-27B Base	Qwen2.5-32B Base	Qwen3-30B-A3B Base	Llama4-Scout Base	GroveMoE Base
Architecture	Dense	Dense	Dense	MoE	MoE	Grove MoE
# Total Params	24B	27B	32B	30B	109B	33B
# Activated Params	24B	27B	32B	3B	17B	3.14B-3.28B
General Tasks						
MMLU	81.65	79.89	83.50	81.58	79.08	<u>82.86</u>
MMLU-Pro	55.67	52.97	59.04	59.58	57.32	<u>59.06</u>
CMMLU	74.85	70.17	88.17	80.63	76.05	<u>86.75</u>
SuperGPQA	30.47	30.15	35.80	<u>36.10</u>	27.54	38.74
BBH	<u>83.46</u>	79.19	84.30	81.58	82.60	82.09
C-Eval	72.81	70.00	86.96	<u>87.82</u>	74.80	87.84
Math & STEM Tasks						
GSM8K	85.90	82.71	<u>90.45</u>	89.39	86.43	90.83
MATH	43.90	49.80	60.42	<u>59.75</u>	51.34	64.82
GPQA-Diamond	39.90	36.36	<u>41.41</u>	39.39	37.54	41.92
Coding Tasks						
HumanEval+	60.98	57.32	78.05	<u>83.54</u>	64.63	85.98
MBPP+	71.16	69.84	<u>73.81</u>	71.96	69.84	76.19
MultiPL-E	27.32	48.20	<u>52.57</u>	61.76	48.53	<u>60.38</u>
CRUX-O	50.38	60.12	<u>67.88</u>	67.20	59.54	70.25

1×10^{-4} to 5×10^{-5} . Consistent with our architectural analysis in Section 4.2, we configure the $\{A_j\}_{j=1}^g$ modules within the Grove MoE layer with group number $g = 64$, intermediate size $h = 128$, and scaling factor $\lambda = 0.05$.

4.4 MID-TRAINING EVALUATION

We compare our GroveMoE-Base models with leading open-source base models, including Mistral-Small-3.1-Base-2503 (Mistral AI, 2025), Gemma3-27B-Base (Gemma et al., 2025), Qwen2.5-32B-Base (Yang et al., 2024), Qwen3-30B-A3B-Base (Yang et al., 2025), and Llama4-Scout-Base (Meta-AI, 2025). All models are evaluated using the same evaluation pipeline and the widely used evaluation settings (Yang et al., 2025) to ensure fair comparison.

As depicted in Table 3, our GroveMoE-Base model, built on the Grove MoE architecture, achieves a strong balance of performance and efficiency. Our Grove MoE architecture enables operation with fewer activated parameters than dense models and achieves greater parameter efficiency than other MoE baselines. GroveMoE-Base excels at complex reasoning, surpassing all competing baselines in Math & STEM and coding benchmarks to achieve the highest average scores. In addition to these advanced reasoning capabilities, GroveMoE-Base is highly competitive in general tasks. It matches the performance of Qwen2.5-32B-Base, a dense model with a similar total parameter count, while maintaining its advantage in activation efficiency.

Notably, GroveMoE-Base is developed based on Qwen3-30B-A3B-Base. As shown in Table 3, the Grove MoE architecture facilitates an efficient expansion of model capacity with only a minor additional computational cost. The Grove MoE architecture allows the model to preserve foundational knowledge while providing dedicated resources to master new, complex skills.

5 POST-TRAINING

5.1 SUPERVISED FINE-TUNING

Following the mid-training stage, the GroveMoE-Base model undergoes supervised fine-tuning (SFT). This stage is crucially dependent on the training data. Given the scarcity and high annotation cost of human-generated data, synthetic data has become increasingly important.

Our SFT dataset is constructed from a seed set of 1–2 million instances, comprising human annotations and open-source materials. This initial dataset is then substantially expanded through a data

Table 4: Comparison among GroveMoE-Inst and other strong open-source non-reasoning baselines. The highest and second-best scores are shown in bold and underlined, respectively.

	Mistral-Small-3.2 Instruct-2506	Gemma3-27B IT	Qwen3-32B Non-Thinking	Qwen3-30B-A3B Non-Thinking	Llama4-Scout	GroveMoE Inst
Architecture	Dense	Dense	Dense	MoE	MoE	Grove MoE
# Total Params	24B	27B	32B	30B	109B	33B
# Activated Params	24B	27B	32B	3B	17B	3.14B-3.28B
General Tasks						
MMLU	80.29	75.97	<u>82.93</u>	80.12	81.88	88.04
MMLU-Pro	68.11	67.10	<u>68.25</u>	63.30	64.92	72.78
CMMU	74.02	65.82	<u>84.63</u>	83.13	76.12	86.66
SuperGPQA	37.53	35.63	<u>43.04</u>	40.50	42.02	47.69
BBH	85.51	<u>85.79</u>	85.45	82.55	77.37	88.42
DROP	86.02	87.81	84.02	86.38	<u>88.26</u>	88.84
C-Eval	72.01	67.31	<u>87.53</u>	85.95	74.69	87.60
AGIEval	58.24	53.63	63.64	<u>65.27</u>	62.31	82.19
Alignment Tasks						
IFEval	82.52	<u>86.14</u>	85.27	84.55	85.57	86.54
Arena-Hard	83.87	<u>89.38</u>	<u>90.49</u>	88.33	73.49	92.01
Math & STEM Tasks						
MATH	84.18	<u>85.82</u>	85.26	84.68	81.46	90.56
MATH-500	86.50	87.80	87.40	<u>88.70</u>	82.60	94.60
Omni-MATH	33.40	33.30	31.80	<u>33.70</u>	25.78	43.50
AIME24	<u>36.88</u>	29.58	27.71	28.33	28.60	54.58
AIME25	<u>28.12</u>	23.12	22.92	21.67	10.00	44.38
GPQA-Diamond	49.94	45.33	53.60	51.71	<u>55.56</u>	61.30
OlympiadBench	<u>61.89</u>	59.85	59.52	60.26	56.11	71.22
Coding & Agent Tasks						
HumanEval+	81.94	78.81	82.93	<u>84.15</u>	79.88	90.24
MBPP+	73.54	73.83	72.75	<u>75.16</u>	70.37	78.31
MultiPL-E	<u>69.49</u>	65.50	68.62	66.04	45.00	74.53
LiveCodeBench v5	25.90	26.75	<u>31.44</u>	28.89	25.45	33.38
LiveCodeBench v6	<u>32.25</u>	30.86	28.57	29.43	32.04	34.60
BFCL v3 (Live)	78.21	75.31	75.09	73.69	45.41	<u>76.11</u>

synthesis pipeline. Our data synthesis process begins by generating novel prompts using methods inspired by Magpie-style approaches (Xu et al., 2024) and OSS-Instruct (Wei et al., 2023). We then apply rejection sampling (Grattafiori et al., 2024) to produce candidate responses using various LLMs (Yang et al., 2025; Google DeepMind, 2025; Hurst et al., 2024). To ensure high data quality, we employ a multi-stage filtering process. Initially, rule-based filters are applied to reasoning-intensive data, such as code, mathematics, and logic problems. Subsequently, all data types undergo a final assessment by an LLM-based evaluator, which uses a detailed rubric to verify the response quality and relevance. This rigorous data curation process yields a robust dataset for SFT.

Moreover, to evaluate the quality of instruction-tuned models, we evaluate LLMs on a series of post-training benchmarks, the details of which are illustrated in Appendix B.2.

5.2 HYPER-PARAMETERS

The post-training stage uses the AdamW optimizer (Loshchilov & Hutter, 2017), with $\beta_1 = 0.9$, $\beta_2 = 0.95$, weight decay of 0.1, and gradient clipping at 1.0. During the post-training stage, we employ a cosine learning rate scheduler with a learning rate of 5×10^{-6} that gradually decays to a minimum of 1×10^{-6} .

5.3 POST-TRAINING EVALUATION

We compare our GroveMoE-Inst with leading open-source LLMs, including Mistral-Small-3.2-Instruct-2506 (Mistral AI, 2025), Gemma3-27B-IT (Gemma et al., 2025), Qwen3-32B (Yang et al., 2025), Qwen3-30B-A3B (Yang et al., 2025), and Llama4-Scout (Meta-AI, 2025).

As shown in Table 4, GroveMoE-Inst establishes excellent performance across a comprehensive set of benchmarks, maintaining high parameter efficiency. In general and alignment tasks, the

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

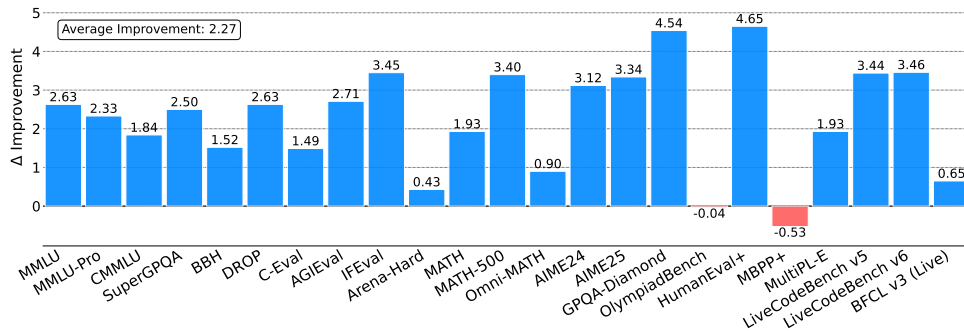


Figure 4: Evaluation results of SFT on various benchmarks. Δ indicates the performance improvement of SFT trained with GroveMoE-Base over Qwen3-30B-A3B-Base.

model consistently outperforms its counterparts, securing the highest scores on all benchmarks. The superiority of our GroveMoE-Inst is particularly pronounced in mathematics and STEM, where GroveMoE-Inst ranks first across all listed benchmarks, highlighting its powerful reasoning capabilities. Furthermore, GroveMoE-Inst demonstrates exceptional performance in coding and agent-based tasks. It surpasses other baselines on most coding & agent benchmarks, which underscores its advanced skills in code generation and problem-solving.

5.4 EFFECTIVENESS OF GROVEMOE-BASE

To evaluate the effectiveness of our GroveMoE-Base model, we applied the same post-training strategy to Qwen3-30B-A3B-Base (Yang et al., 2025) for a direct comparison. [We also provide additional ablation study in Appendix C.2 for fair comparison under the same training corpus.](#)

As shown in Figure 4, the instruction-tuned model derived from GroveMoE-Base consistently outperforms its Qwen3-30B-A3B-Base counterpart across the vast majority of tasks, highlighting its strong potential as a foundation model. Specifically, the GroveMoE-Base model achieves higher scores on nearly all general and alignment benchmarks. The advantage is further pronounced in specialized domains, where it significantly outperforms the Qwen model on most mathematics and STEM benchmarks. Furthermore, its superiority extends to code generation and agent-based tasks, where it secures stronger results on key benchmarks.

In summary, these results demonstrate that GroveMoE-Base is a more powerful foundation model. Its larger model capacity enables fine-tuned derivatives to achieve superior performance across a wide spectrum of domains, including general knowledge, mathematics, and coding.

6 CONCLUSION

This paper introduces GroveMoE models, efficient and open-source LLMs built upon the Grove MoE architecture, which incorporates a novel mechanism for dynamic computation allocation. The Grove MoE architecture improves computational efficiency by dividing experts into groups, each with an adjugate expert. This design ensures that shared computations are performed only once per group, even when multiple experts are activated. GroveMoE-Base and GroveMoE-Inst are 33B-parameter models developed based on the Qwen3-30B-A3B-Base model using our Grove MoE architecture during the mid-training and post-training stage. It dynamically activates 3.14–3.28B parameters per token. Empirical evaluations demonstrate that our GroveMoE models, including Base and Inst models, achieve performance comparable to SOTA open-source models of similar or even larger sizes, thereby validating the effectiveness of the Grove MoE architecture.

DECLARATION OF LLM USAGE

The usage of LLMs is strictly limited to aid and polish the paper writing.

REFERENCES

- 486
487
488 AIME. AIME Problems and Solutions, 2025. URL https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions.
489
- 490 Cenk Baykal, Dylan Cutler, Nishanth Dikkala, Nikhil Ghosh, Rina Panigrahy, and Xin Wang. Alternating Updates for Efficient Transformers. In *Annual Conference on Neural Information Processing Systems (NIPS)*, 2023.
491
492
493
- 494 ByteDance-Seed. Seed-OSS Open-Source Models. <https://github.com/ByteDance-Seed/seed-oss>, 2025.
495
- 496 Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, et al. MultiPLE: A Scalable and Polyglot Approach to Benchmarking Neural Code Generation. *IEEE Transactions on Software Engineering*, 49(7):3675–3691, 2023.
497
498
499
- 500 Mouxiang Chen, Binyuan Hui, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Jianling Sun, Junyang Lin, and Zhongxin Liu. Parallel Scaling Law for Language Models. *arXiv preprint arXiv:2505.10475*, 2025.
501
502
503
- 504 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training Verifiers to Solve Math Word Problems. *arXiv preprint arXiv:2110.14168*, 2021.
505
506
507
- 508 Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, King Zhu, Minghao Liu, Yiming Liang, Xiaolong Jin, Zhenlin Wei, et al. SuperGPQA: Scaling llm evaluation across 285 graduate disciplines. *arXiv preprint arXiv:2502.14739*, 2025.
509
510
- 511 Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. DROP: A Reading Comprehension Benchmark Requiring Discrete Reasoning over Paragraphs. *arXiv preprint arXiv:1903.00161*, 2019.
512
513
514
- 515 Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li, Chenghao Ma, Liang Chen, Runxin Xu, et al. Omni-MATH: A Universal Olympiad Level Mathematic Benchmark for Large Language Models. *arXiv preprint arXiv:2410.07985*, 2024.
516
517
- 518 Gemma, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 Technical Report. *arXiv preprint arXiv:2503.19786*, 2025.
519
520
521
- 522 Google DeepMind. Gemini2.5 Pro. <https://deepmind.google/technologies/gemini/pro/>, 2025.
523
- 524 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The Llama 3 Herd of Models. *arXiv preprint arXiv:2407.21783*, 2024.
525
526
527
- 528 Peter Greenhalgh. big.LITTLE Processing with ARL Cortex-A15 & Cortex-A7. *ARM White paper*, 17, 2011.
529
- 530 Alex Gu, Baptiste Rozière, Hugh Leather, Armando Solar-Lezama, Gabriel Synnaeve, and Sida I Wang. CruxEval: A Benchmark for Code Reasoning, Understanding and Execution. *arXiv preprint arXiv:2401.03065*, 2024.
531
532
533
- 534 Yongxin Guo, Zhenglin Cheng, Xiaoying Tang, Zhaopeng Tu, and Tao Lin. Dynamic Mixture of Experts: An Auto-Tuning Approach for Efficient Transformer Models. *arXiv preprint arXiv:2405.14297*, 2024.
535
536
- 537 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. OlympiadBench: A Challenging Benchmark for Promoting AGI with Olympiad-Level Bilingual Multimodal Scientific Problems. *arXiv preprint arXiv:2402.14008*, 2024.
538
539

- 540 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
541 Jacob Steinhardt. Measuring Massive Multitask Language Understanding. *arXiv preprint*
542 *arXiv:2009.03300*, 2020.
- 543
544 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
545 and Jacob Steinhardt. Measuring Mathematical Problem Solving with the Math Dataset. *arXiv*
546 *preprint arXiv:2103.03874*, 2021.
- 547
548 Quzhe Huang, Zhenwei An, Nan Zhuang, Mingxu Tao, Chen Zhang, Yang Jin, Kun Xu, Liwei
549 Chen, Songfang Huang, and Yansong Feng. Harder Tasks Need More Experts: Dynamic Routing
550 in MoE Models. *arXiv preprint arXiv:2403.07652*, 2024.
- 551
552 Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu,
553 Chuancheng Lv, Yikai Zhang, Yao Fu, et al. C-Eval: A Multi-Level Multi-Discipline Chinese
554 Evaluation Suite for Foundation Models. *Advances in Neural Information Processing Systems*,
36:62991–63010, 2023.
- 555
556 Bi Huo, Bin Tu, Cheng Qin, Da Zheng, Debing Zhang, Dongjie Zhang, En Li, Fu Guo, Jian Yao,
557 Jie Lou, et al. dots.llm1 Technical Report. *arXiv preprint arXiv:2506.05767*, 2025.
- 558
559 Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Os-
560 trow, Akila Welihinda, Alan Hayes, Alec Radford, et al. GPT-4o System Card. *arXiv preprint*
arXiv:2410.21276, 2024.
- 561
562 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando
563 Solar-Lezama, Koushik Sen, and Ion Stoica. LiveCodeBench: Holistic and Contamination Free
564 Evaluation of Large Language Models for Code. *arXiv preprint arXiv:2403.07974*, 2024.
- 565
566 Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bam-
567 ford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al.
Mixtral of Experts. *arXiv preprint arXiv:2401.04088*, 2024.
- 568
569 Peng Jin, Bo Zhu, Li Yuan, and Shuicheng Yan. MoE++: Accelerating Mixture-of-Experts Methods
570 with Zero-Computation Experts. *arXiv preprint arXiv:2410.07348*, 2024.
- 571
572 Aran Komatsuzaki, Joan Puigcerver, James Lee-Thorp, Carlos Riquelme Ruiz, Basil Mustafa,
573 Joshua Ainslie, Yi Tay, Mostafa Dehghani, and Neil Houlsby. Sparse Upcycling: Training
574 mixture-of-experts from dense checkpoints. In *International Conference on Learning Repre-*
sentations (ICLR), 2023.
- 575
576 Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy
577 Baldwin. CMMLU: Measuring Massive Multitask Language Understanding in Chinese. *arXiv*
preprint arXiv:2306.09212, 2023.
- 578
579 Tianle Li, Wei-Lin Chiang, Evan Frick, Lisa Dunlap, Tianhao Wu, Banghua Zhu, Joseph E Gon-
580 zalez, and Ion Stoica. From Crowdsourced Data to High-Quality Benchmarks: Arena-Hard and
581 Benchbuilder Pipeline. *arXiv preprint arXiv:2406.11939*, 2024.
- 582
583 Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan
584 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s Verify Step by Step. In *The Twelfth*
International Conference on Learning Representations, 2023.
- 585
586 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,
587 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. DeepSeek-V3 Technical Report. *arXiv preprint*
588 *arXiv:2412.19437*, 2024.
- 589
590 Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is Your Code Generated by
591 ChatGPT Really Correct? Rigorous Evaluation of Large Language Models for Code Generation.
Advances in Neural Information Processing Systems, 36:21558–21572, 2023.
- 592
593 Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. *arXiv preprint*
arXiv:1711.05101, 2017.

- 594 Meta-AI. The Llama 4 herd: The beginning of a new era of natively multimodal AI innovation, 2025.
595 URL <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>.
- 596
- 597 Mistral AI. Mistral-Small-3.1. <https://mistral.ai/news/mistral-small-3-1>, 2025.
- 598
- 599 Taishi Nakamura, Takuya Akiba, Kazuki Fujii, Yusuke Oda, Rio Yokota, and Jun Suzuki. Drop-
600 Upcycling: Training Sparse Mixture of Experts with Partial Re-Initialization. *arXiv preprint*
601 *arXiv:2502.19261*, 2025.
- 602 David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien
603 Dirani, Julian Michael, and Samuel R Bowman. GPQA: A Graduate-Level Google-Proof Q&Q
604 Benchmark. In *First Conference on Language Modeling*, 2024.
- 605 sgl-project. SGLang. <https://github.com/sgl-project/sglang>, 2025.
- 606
- 607 Jianlin Su. MoE Travels 3, 2025. URL <https://kexue.fm/archives/10757>.
- 608
- 609 Manxi Sun, Wei Liu, Jian Luan, Pengzhi Gao, and Bin Wang. Mixture of Diverse Size Experts.
610 *arXiv preprint arXiv:2409.12210*, 2024.
- 611 Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung,
612 Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. Challenging Big-Bench Tasks
613 and Whether Chain-of-Thought Can Solve Them. *arXiv preprint arXiv:2210.09261*, 2022.
- 614 An Wang, Xingwu Sun, Ruobing Xie, Shuaipeng Li, Jiaqi Zhu, Zhen Yang, Pinxue Zhao, Weidong
615 Han, Zhanhui Kang, Di Wang, et al. HMoE: Heterogeneous Mixture of Experts for Language
616 Modeling. In *Empirical Methods in Natural Language Processing (EMNLP)*, pp. 21954–21968,
617 2025.
- 618
- 619 Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming
620 Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. MMLU-Pro: A More Robust and Challeng-
621 ing Multi-Task Language Understanding Benchmark. In *The Thirty-eight Conference on Neural*
622 *Information Processing Systems Datasets and Benchmarks Track*, 2024a.
- 623 Ziteng Wang, Jun Zhu, and Jianfei Chen. ReMoE: Fully Differentiable Mixture-of-Experts with
624 ReLU Routing. *arXiv preprint arXiv:2412.14711*, 2024b.
- 625
- 626 Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Empowering
627 Code Generation with OSS-Instruct. *arXiv preprint arXiv:2312.02120*, 2023.
- 628
- 629 Haoyuan Wu, Haisheng Zheng, Zhuolun He, and Bei Yu. Parameter-Efficient Sparsity Crafting from
630 Dense to Mixture-of-Experts for Instruction Tuning on General Tasks. In *Empirical Methods in*
631 *Natural Language Processing (EMNLP)*, 2024.
- 632
- 633 Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and
634 Bill Yuchen Lin. Magpie: Alignment Data Synthesis from Scratch by Prompting Aligned LLMs
635 with Nothing. *arXiv preprint arXiv:2406.08464*, 2024.
- 636
- 637 Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica, and
638 Joseph E. Gonzalez. Berkeley Function Calling Leaderboard. https://gorilla.cs.berkeley.edu/blogs/8_berkeley_function_calling_leaderboard.html,
639 2024.
- 640
- 641 An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,
642 Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2.5 Technical Report. *arXiv preprint*
643 *arXiv:2412.15115*, 2024.
- 644
- 645 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
646 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 Technical Report. *arXiv preprint*
647 *arXiv:2505.09388*, 2025.
- 648
- 649 Zihao Zeng, Yibo Miao, Hongcheng Gao, Hao Zhang, and Zhijie Deng. AdaMoE: Token-
650 Adaptive Routing with Null Experts for Mixture-of-Experts Language Models. *arXiv preprint*
651 *arXiv:2406.13233*, 2024.

648 Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied,
649 Weizhu Chen, and Nan Duan. AGIEval: A Human-Centric Benchmark for Evaluating Foundation
650 Models. *arXiv preprint arXiv:2304.06364*, 2023.

651
652 Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny
653 Zhou, and Le Hou. Instruction-Following Evaluation for Large Language Models. *arXiv preprint*
654 *arXiv:2311.07911*, 2023.

655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A MOE ARCHITECTURE

A.1 DYNAMIC ACTIVATION

Existing dynamic activation methods in MoE models aim to solve the limitations of fixed top- k routing. However, current approaches either indirectly vary expert counts using blank experts (Jin et al., 2024; Zeng et al., 2024) or employ complex routing strategies (Guo et al., 2024; Wang et al., 2024b) that suffer from uncontrollable computational overhead and incompatibility with standard top- k mechanisms. In contrast, as shown in Table 5, our GroveMoE intrinsically achieves dynamic activation through expert grouping. This design ensures that both the number of active experts and the computational cost are controllable. Furthermore, it requires no router modifications, making it highly compatible and easy to implement.

Table 5: Comparison among MoE Architectures with Dynamic Activation.

Scheme	Dynamism	Computational Control	Compatibility	Parameter Efficiency
Fixed Top-k (Jin et al., 2024; Zeng et al., 2024)	✗ Rigid	✓ Strict	✓ High	✗ Low
Dynamic Top-k (Guo et al., 2024; Wang et al., 2024b)	✓ Flexible	✗ Explicit Constraints	✗ Redesign of the Router	✗ Low
GroveMoE (Ours)	✓ Flexible	✓ Implicit Grouping	✓ Compatible with Top- k	✓ Shared Adjudication

A.2 HETEROGENEOUS EXPERTS

The MoE layer with heterogeneous experts proposes that experts should vary in size to possess diverse capacities (Wang et al., 2025; Sun et al., 2024). This heterogeneity allows the model to route input tokens to the expert with appropriate capability, thereby more effectively handling the varying complexity of the data and enabling greater expert specialization. In this research, our utilization of adjugate experts with varying sizes serves as a deliberate design choice to control the granularity of the activated parameter counts. This structure ensures that the step size for dynamic activation is neither excessively large nor too small, enabling finer control over the activated model capacity.

B EVALUATION BENCHMARKS

B.1 EVALUATION BENCHMARKS FOR MID-TRAINING

Our comprehensive evaluation of the base models assesses five core capabilities: general knowledge, scientific knowledge, reasoning, mathematics, and coding. The evaluation is conducted using 13 distinct benchmarks:

- **General Tasks:** MMLU (Hendrycks et al., 2020)(5-shot), MMLU-Pro (Wang et al., 2024a)(5-shot, CoT), CMMLU (Li et al., 2023)(5-shot), SuperGPQA (Du et al., 2025)(5-shot), C-Eval (Huang et al., 2023)(5-shot), and BBH (Suzgun et al., 2022)(3-shot, CoT).
- **Math & STEM Tasks:** GSM8K (Cobbe et al., 2021)(4-shot, CoT), MATH (Hendrycks et al., 2021)(4-shot, CoT), and GPQA-Diamond (Rein et al., 2024)(5-shot).
- **Coding Tasks:** HumanEval+ (Liu et al., 2023)(0-shot), MBPP+ (Liu et al., 2023)(0-shot), MultiPL-E (Cassano et al., 2023)(0-shot)(Python, C++, Java, PHP, TypeScript, C#, Bash, JavaScript), and CRUX-O (Gu et al., 2024)(1-shot, CoT).

B.2 EVALUATION BENCHMARKS FOR POST-TRAINING

To evaluate the quality of instruction-tuned models, we evaluate LLMs on a series of post-training benchmarks. The post-training benchmarks can be categorized into several dimensions:

- **General Tasks:** For general language understanding tasks, we utilize various benchmarks including MMLU (Hendrycks et al., 2020), MMLU-Pro (Wang et al., 2024a), CMMLU (Li et al., 2023), SuperGPQA (Du et al., 2025), BBH (Suzgun et al., 2022), DROP (Dua et al., 2019), C-Eval (Huang et al., 2023), and AGIEval (Zhong et al., 2023).

Table 6: Ablation studies on different routing strategies.

	-	-	$g = 64; h = 128$ (Ours)	$g = 128; h = 128$	$g = 64; h = 128$
Architecture	MoE	MoE	Grove MoE	-	-
Routing Strategy	Single Router	Single Router	Single Router	Single Router	Coupled Router
# Total Params	30B	30B	33B	36B	33B
# Activated Params	3B	3B	3.14B-3.28B	3.28B	3.28B
# Avg. Activation	3B	3B	3.26B	3.28B	3.28B
# Training Tokens	0B	50B	50B	50B	50B
MMLU	81.58	82.56	82.62	81.89	82.06
CMMLU	80.63	86.54	86.55	86.75	86.43
SuperGPQA	36.10	35.93	36.32	36.09	36.50
GSM8K	89.39	89.54	90.83	91.05	90.83
MATH	59.75	65.90	65.86	66.03	65.94
GPQA-Diamond	39.39	38.38	43.94	40.91	42.42
HumanEval+	83.54	83.54	84.76	84.15	85.37
MBPP+	71.96	74.34	75.13	75.40	75.13
Average	67.79	69.59	70.75	70.28	70.58

- Alignment Tasks:** To evaluate how well the model aligns with human preferences, we employ a suite of specialized benchmarks. For instruction-following performance, we report the average prompt-level and instruction-level strict accuracy of IFEval (Zhou et al., 2023). To assess alignment with human preferences on general topics, we utilize Arena-Hard (Li et al., 2024).
- Math & STEM Tasks:** For evaluating mathematical reasoning skills, we employ MATH (Hendrycks et al., 2021), MATH-500 (Lightman et al., 2023), Omni-MATH (Gao et al., 2024), AIME24 (AIME, 2025), and AIME25 (AIME, 2025). For STEM tasks, we utilize GPQA-Diamond (Rein et al., 2024) and OlympiadBench (He et al., 2024) as evaluation benchmarks. For AIME problems, we sample 16 times for each question and take the average accuracy as the final score. For GPQA-Diamond, we sample 8 times for each query and report the average accuracy.
- Coding & Agent Tasks:** To test the LLM’s proficiency in coding and agent-based tasks, we use HumanEval+ (Liu et al., 2023), MBPP+ (Liu et al., 2023), MultiPL-E (Cassano et al., 2023), LiveCodeBench (Jain et al., 2024) (v5, 2024.10-2025.02 and v6, 2025.02-2025.05), and BFCL v3 (Live) (Yan et al., 2024). For BFCL v3 (Live), all models are evaluated using the prompt format.

Notably, we configure the maximum output length to 16K to avoid overly lengthy output for non-reasoning LLMs during the evaluation process (ByteDance-Seed, 2025; Yang et al., 2025).

C ABLATION STUDIES

C.1 ROUTING STRATEGY

To further investigate the effectiveness of our routing strategy, we perform a series of ablation studies. As presented in Table 6, our GroveMoE architecture with dynamic activation achieves performance comparable to architectures employing fixed expert activation. Specifically, the GroveMoE architecture maintains competitive performance while ensuring efficiency through adaptive computation. This finding demonstrates the efficacy and necessity of the adaptive routing mechanism in mitigating the performance disparity while effectively controlling computational cost. The competitive performance of GroveMoE, even when compared to an architecture that activates a fixed number of adjudicate experts, is noteworthy. This outcome suggests that our GroveMoE successfully learns to discern token complexity, enabling dynamic computation.

C.2 ENTIRE TRAINING FLOW

To ensure a fair comparison under the same training conditions, we conducted an additional ablation study, the results of which are presented in Table 7. We perform this study using the Qwen3-30B-A3B-Base model, undergoing mid-training for 50B tokens, and applying the same post-training

Table 7: Ablation studies on the entire training flow. The highest score is shown in bold.

Models	Qwen3-30B-A3B	GroveMoE (Ours)
# Total Params	30B	33B
# Activated Params	3B	3.14B-3.28B
# Avg. Activation	3B	3.26B
# Training Tokens	50B	50B
# SFT Epochs	1	1
MMLU	82.23	83.50
MMLU-Pro	67.21	68.91
SuperGPQA	42.02	45.39
BBH	84.14	85.15
CEval	84.23	83.93
IFEval	81.54	85.74
GSM8K	90.45	93.03
MATH-500	86.40	88.30
GPQA-Diamond	56.19	58.21
HumanEval+	85.06	85.67
MBPP+	75.93	76.98
MultiPLE	68.27	69.47
LiveCodeBench v5	26.80	27.54
Average	71.57	72.45

setup for one epoch. As shown in Table 7, our GroveMoE consistently outperforms the Qwen3-30B-A3B baseline across most metrics, which further affirms its architectural benefits.

D GROUPING STRATEGY

The expert grouping strategy represents a fascinating and potentially important research avenue. As detailed in Equation (2) and Equation (3), the grouping mechanism within our GroveMoE architecture utilizes a simple partition. Although a grouping strategy predicated on co-activation patterns or functional similarity could yield substantial benefits, we deliberately avoided an in-depth, data-driven analysis during training. This choice mitigates potential methodological issues, such as introducing bias or data leakage when using a validation set to determine the grouping. The significant gains in performance and efficiency achieved, despite employing this relatively simple grouping method, underscore the robustness of the core GroveMoE architecture.

E DEPLOYMENT OF GROVEMOE-INST

In our native PyTorch implementation, the Grove MoE architecture achieves the claimed theoretical speed. However, in our SGLang (sgl-project, 2025) implementation, the inference speed of GroveMoE-Inst is approximately 30% slower than the Qwen3-30B-A3B (Yang et al., 2025) baseline, which is an overhead that significantly exceeds the theoretical maximum 10% increase in activated parameters. This discrepancy arises because our implementation uses the generic MoE kernel (namely fused-moe) of SGLang for accessibility, which requires two separate kernel calls per GroveMoE layer. A customized and unified kernel designed to process both expert types in one operation would mitigate this latency and align performance more closely with theoretical expectations. The development of such a kernel is our key priority for future efficiency improvements.

LIMITATIONS

Although our Grove MoE architecture provides a solid foundation, two primary limitations constrain its current potential and guide our future research. The first limitation stems from a scarcity of long-CoT data within the mid-training corpus. This data deficiency curtails the model’s capacity for advanced reasoning, creating a capability gap compared to instruction-tuned LLMs that possess stronger foundational reasoning abilities, such as Qwen3-30B-A3B-2507 (Yang et al., 2025), etc. The second limitation is the exclusive reliance on rejection sampling for model refinement, without the integration of RL techniques. While rejection sampling has been effective, we anticipate that

864 incorporating RL methods will significantly enhance the model’s overall capabilities. This remains
865 a key objective for future development.
866

867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917