

Router-Tuning: A Simple and Effective Approach for Dynamic Depth

Anonymous ACL submission

Abstract

The Mixture of Depths (MoD) was introduced to improve computational efficiency by dynamically skipping less important layers, reducing redundant computation while maintaining model capacity. Despite its promise, existing MoD approaches remain under-explored and face two main challenges: (1) *high training costs due to the need to train the entire model along with the routers that determine which layers to skip*, and (2) *performance degradation when important layers are bypassed*. In response to the first issue, we propose Router-Tuning, which fine-tunes only the routers on a small dataset, drastically reducing the computational overhead associated with full model training. For the second challenge, we investigate Router-Tuning across different architectures and granularities, demonstrating its effectiveness on Attention layers and MoE layers. This method preserves the model’s performance while significantly enhancing computational and memory efficiency. Extensive experiments demonstrate that our approach delivers competitive results while dramatically improving the computation efficiency, e.g., 21% speedup and only a 0.2% performance drop. The code will be released upon acceptance.

1 Introduction

Large Language Models (LLMs) have shown promising performance across various domains (OpenAI et al., 2024; Team, 2024; DeepSeek-AI et al., 2024). However, the continuous increase in model size leads to substantial computational costs in real-world applications, making computation reduction a critical research focus for improving LLM efficiency (Sun et al., 2024; Lin et al., 2024). A promising approach to this challenge is the Mixture of Depths (MoD) (Raposo et al., 2024), which dynamically allocates computational resources for specific inputs. Instead of uniformly applying all layers to every input, MoD selectively activates

only a subset of the model’s layers, skipping those deemed less important. This targeted activation significantly reduces computational overhead while maintaining performance.

Despite its potential, current MoD methods are still underexplored and face several critical challenges. On the one hand, the involvement of additional router networks, which decide which layers to skip, often requires extra extensive training: Raposo et al. (2024) train the entire model from scratch while Tan et al. (2024) performs costly continual training. This creates a significant barrier to efficiently integrating MoD with existing LLMs. Furthermore, most prior MoD implementations (Raposo et al., 2024; Tan et al., 2024) have been applied to transformer blocks and MLP layers, which are sensitive to skipping. As a result, omitting important components often leads to significant performance degradation (He et al., 2024b).

These challenges prompt us to reflect on the two key questions: (1) *How can we implement dynamic depth to improve efficiency without incurring excessive training costs?* (2) *How can we preserve model performance in the presence of dynamic depth?*

To tackle the first challenge, we introduce *Router-Tuning*, a novel method that fine-tunes only the router network without updating the backbone model’s parameters. As each router network is a lightweight, single-layer projector that accounts for less than 0.01% of the total parameters, the training overhead is minimal and even significantly lower than that of parameter-efficient finetuning methods (Houlsby et al., 2019a; He et al., 2021) like LoRA (Hu et al., 2022). Router-tuning requires only a small-scale dataset and fewer training steps, eliminating the need for large-scale pretraining or extensive continual training. Meanwhile, by freezing the backbone, Router-Tuning contributes to mitigating catastrophic forgetting and better retaining the original model performance (Houlsby et al., 2019b; Liu et al., 2024; Qiao and Mahdavi, 2024). These

properties make Router-Tuning a highly efficient and scalable solution for dynamic adaptation.

To address the second challenge, we conduct a comprehensive investigation of target modules (e.g., Block, MLP, Attention, MoE), and various granularities (e.g., token and sequence). For dense transformer architectures, we propose *Attention with Dynamic Depths*, which selectively applies dynamic depth to attention layers. By focusing on attention layers known to exhibit high redundancy (He et al., 2024b), Router-Tuning not only preserves model accuracy but also alleviates computational and memory bottlenecks. In the case of Mixture-of-Experts (MoE) layers (Shazeer et al., 2017; Fedus et al., 2022), where efficiency is often hindered by the computational cost of activating multiple expert networks, we apply Router-Tuning at the expert level to enhance overall efficiency.

Through extensive experiments, we demonstrate the effectiveness of our approach across multiple open-source language models, including Llama (Touvron et al., 2023), Mistral (Jiang et al., 2023), Qwen (Bai et al., 2023), Deepseek-MoE (Dai et al., 2024), and OLMoE (Muennighoff et al., 2024). Router-Tuning requires less than half an hour on an Nvidia RTX A6000, making it 1000 times faster than DLO (Tan et al., 2024). Router-Tuning maintains a high percentage of the original model’s performance while significantly reducing memory usage and accelerating inference, achieving, for example, a 21% inference speedup with only a 0.2% performance degradation. Furthermore, Router-Tuning can be seamlessly integrated with LoRA fine-tuning, further enhancing both efficiency and performance.

In short, our contributions are as follows:

- We introduce *Router-Tuning*, a lightweight method that fine-tunes only the router using a small dataset, effectively addressing the high computational cost of training the entire model with routers.
- We systematically investigate routing scopes, deployment granularities, and model architectures, demonstrating the effectiveness of Router-Tuning on Attention and MoE layers.
- Through comprehensive experiments, Router-Tuning achieves competitive performance while delivering substantial improvements in training and inference efficiency.

2 Related Work

Layer Redundancy While increasing the depth of large language models has demonstrated promising performance across a wide range of tasks (OpenAI et al., 2024; Team, 2024), it also introduces layer redundancy (Gromov et al., 2024; He et al., 2024b), posing efficiency challenges. To address this issue, several approaches have been proposed to reduce model depth (Men et al., 2024) while maintaining comparable performance. Surprisingly, removing redundant layers has been shown to preserve performance while significantly reducing memory and computational costs (Gromov et al., 2024; He et al., 2024a,b). Specifically, Gromov et al. (2024) suggest dropping continuous Transformer blocks, and He et al. (2024b) propose fine-grained layer dropping to further improve the effectiveness of layer reduction. However, these static techniques fail to account for the varying complexity of different input sequences, where excessive layer removal can significantly degrade performance on more complex tasks. Instead of statically removing unimportant layers, our approach focuses on dynamically skipping these layers based on the specific inputs.

Dynamic Depth Dynamic Depth, which allocates different layers based on the specific input, is an effective technique for accelerating inference while preserving performance (Han et al., 2021, 2022). Recent works primarily implement dynamic depth through two key methods: Early-Exit (Bae et al., 2023; Elhoushi et al., 2024) and Mixture of Depths (MoD) (Raposo et al., 2024). Early-exit strategies terminate computation in later layers once sufficient confidence is achieved, effectively reducing redundant computations. In contrast, MoD offers greater flexibility by dynamically skipping less critical layers, enhancing adaptability and representational capacity. Despite their advantages, both Early-Exit and MoD often involve significant training overhead. For instance, LayerSkip (Elhoushi et al., 2024) and MoD (Raposo et al., 2024) require training models from scratch or extensive continual training, while Tan et al. (Tan et al., 2024) extends pre-trained model training over long schedules to achieve optimal performance. To overcome these limitations, we propose Router-Tuning, an efficient approach to dynamic layer skipping that requires minimal additional offline training, providing a more cost-effective solution.

3 Methodology

In this section, we first review the challenges associated with deploying Mixture of Depths and then introduce Router-Tuning, addressing the implementation of Mixture of Depths from both design and training perspectives.

3.1 Motivation

The Mixture of Depths (MoD) framework (Raposo et al., 2024), which dynamically adjusts layer depth based on input complexity to enhance computational efficiency, was originally designed for integration during the pretraining phase, where transformer models are trained from scratch with MoD-enabled layers. More recently, Tan et al. (2024) applied MoD to pretrained Llama models (Touvron et al., 2023) through continual training. While these approaches have demonstrated promising results, *training with MoD remains computationally expensive and time-consuming, posing challenges for scalability and real-world deployment*. A more efficient alternative is to apply MoD directly to existing pretrained models, followed by lightweight fine-tuning of a subset of parameters (Houlsby et al., 2019b; Hu et al., 2022), significantly reducing both computational costs and training time.

On the other hand, MoD has typically been implemented at the transformer block level. However, *skipping entire transformer blocks has shown to be suboptimal to maintain the performance*. Inspired by He et al. (2024b), we recognize that each transformer block contains layers of varying importance. Aggressively skipping entire blocks risks omitting critical layers, potentially degrading performance. Instead, skipping fine-grained layers offers a more effective strategy for preserving model accuracy. Moreover, unlike blocks that generally share the same architecture, individual layers impose different computational costs. For instance, in dense transformer models, attention layers are particularly expensive, with computational complexity scaling quadratically with sequence length and additional memory needed for KV cache storage. In contrast, in MoE models, MLP layers hold the majority of the parameters, leading to substantial communication and computation overhead.

Building on these insights, we propose Router-Tuning, a cost-effective formulation of MoD that achieves a favorable trade-off between performance and computational costs.

3.2 Router-Tuning for Dynamic Depth

In this part, we propose Router-Tuning to address the challenges outlined in Section 3.1. As illustrated in Figure 1, Router-Tuning incorporates an additional trainable router that determines whether to skip the layer. Specifically, Router-Tuning can be deployed in two levels: (1) *token-level*, where layers are dynamically skipped for individual tokens, and (2) *sequence-level*, where layers are dynamically skipped for the entire sequence. Given an input $\mathbf{x} \in \mathbb{R}^{L \times d}$, the router first computes an importance score for the input:

$$\mathbf{R}(\mathbf{x})_i = \begin{cases} \text{GATE}(\mathbf{x}_i), & \text{Token-level} \\ \text{GATE}(\frac{1}{L} \sum_{i=1}^L \mathbf{x}_i), & \text{Sequence-level} \end{cases}, \quad (1)$$

where \mathbf{R} is a scoring router that assesses the importance score of the input, GATE is the gating function $\text{GATE}(\mathbf{x}) = \text{Sigmoid}(\mathbf{W}\mathbf{x})$. Based on the computed importance scores, we further apply a binarized mask \mathbf{M} to determine whether to skip a token or an entire sequence:

$$\mathbf{M} = \begin{cases} 1, & \text{if } \mathbf{R}(\mathbf{x}) \geq \tau \\ 0, & \text{otherwise} \end{cases}, \quad (2)$$

where τ is the threshold. The score is set to zero for skipped inputs and one for retained inputs, ensuring stable outputs (Tan et al., 2024).

To enable a differentiable and trainable binary decision process, we utilize the straight-through estimator (STE) (Bengio et al., 2013), which allows gradients to propagate through the binary selection mechanism via $\frac{\partial \mathbf{M}}{\partial \mathbf{R}} = 1$. The final output of MoD is then computed as follows:

$$\mathbf{y} = \mathbf{M} \odot \mathbf{F}(\mathbf{x}) + \mathbf{x}, \quad (3)$$

where \mathbf{F} denotes a given layer and \mathbf{y} is the output. This formulation ensures that the router is fully trainable through the gradient calculations:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{W}} = \frac{\partial \mathbf{y}}{\partial \mathbf{M}} \frac{\partial \mathbf{M}}{\partial \mathbf{R}} \frac{\partial \mathbf{R}}{\partial \mathbf{W}}. \quad (4)$$

During inference, without the need for gradient calculations, we further enhance computational efficiency by completely bypassing computations for skipped inputs:

$$\mathbf{y} = \begin{cases} \mathbf{F}(\mathbf{x}) + \mathbf{x}, & \text{if } \mathbf{R}(\mathbf{x}) \geq \tau \\ \mathbf{x}, & \text{otherwise} \end{cases}. \quad (5)$$

This dynamic routing mechanism ensures that computation is performed only when necessary, thereby enhancing the computational efficiency.

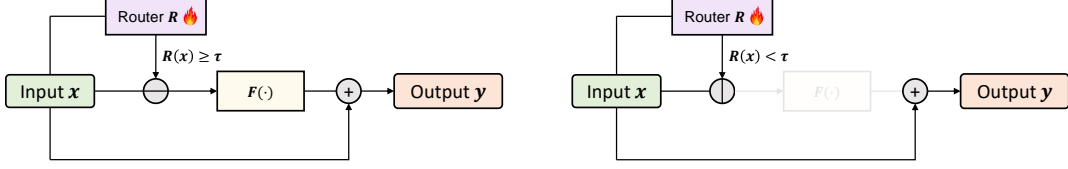


Figure 1: **Overview of Router-Tuning.** Router-Tuning involves a trainable router to determine whether a given layer $F(\cdot)$ (e.g., Attention and MLP) would be skipped. Inputs with routing scores lower $R(x)$ than the threshold τ are skipped, and only the router R is trainable in the whole model.

3.3 Extension to Mixture of Experts

Mixture of Experts (MoE) employs sparse activation, dynamically selecting expert networks for each input, which delivers promising performance in various tasks (Jiang et al., 2024; Dai et al., 2024; Muennighoff et al., 2024). However, MoE also exhibits significant redundancy, allowing certain experts or layers to be skipped with minimal impact on performance (Lu et al., 2024; He et al., 2024a). Building on this, we extend Router-Tuning to MoE layers by implementing dynamic skipping within each expert:

$$\hat{E}_i(x) = \begin{cases} E_i(x), & \text{if } R(x) \geq \tau \\ 0, & \text{otherwise} \end{cases}, \quad (6)$$

where E_i denotes the i -th expert and $\hat{E}_i(x)$ is the corresponding output denotes the corresponding output, bypassing the skipped tokens. Given a collection of n experts, $\{E_1, E_2, \dots, E_n\}$, the overall output of the MoE layer is as follows:

$$\mathcal{K} = \text{TopK}(\text{Softmax}(G(x)), k), \quad (7)$$

$$y = \sum_{i \in \mathcal{K}} G(x)_i \cdot \hat{E}_i(x), \quad (8)$$

where \mathcal{K} denotes the indices of the top- k selected experts, and $G(x)_i$ represents the selection score for the i -th expert. By dynamically skipping experts within each layer, Router-Tuning significantly reduces computation costs.

3.4 Training Objectives

Given the computationally intensive nature of training entire LLMs and the constraints of real-world computational resources, our goal is to implement dynamic depth while minimizing both computational costs and time overhead. To achieve this, we focus exclusively on fine-tuning the routers, as illustrated in Figure 1, thereby eliminating the need for costly full-model training.

Specifically, we optimize two training objectives: improving task-specific performance and lowering

MoD capacity (the proportion of non-skipped inputs). On the one hand, Router-Tuning is designed to maintain the performance of the original model, which we enforce using the loss term $\mathcal{L}_{\text{task}}$ during fine-tuning. On the other hand, the model is encouraged to skip more tokens or sequences (i.e., reduce MoD capacity) to enhance efficiency. To achieve this, we introduce another loss term \mathcal{L}_{MoD} , which drives the model to reduce MoD capacity to a desired target sparsity level s , thereby lowering computational costs and accelerating inference. The overall training objective is as follows:

$$\mathcal{L} = \mathcal{L}_{\text{task}} + \lambda \cdot \mathcal{L}_{\text{MoD}}, \quad (9)$$

$$\mathcal{L}_{\text{MoD}} = \text{ReLU}(\|M\|_0 - s), \quad (10)$$

where \mathcal{L} represents the standard loss function (e.g., cross-entropy), while \mathcal{L}_{MoD} is an l_0 -norm regularization term that reduces MoD capacity. The coefficient λ acts as a scaling factor to balance the trade-off between task performance and efficiency.

4 Main Results

In this section, we evaluate the effectiveness of Router-Tuning on transformer architectures with the deployment details in Appendix A.

4.1 Performance of Router-Tuning

Router-Tuning achieves superior performance on Attention layers We first compare deploying Router-Tuning to different modules, e.g., Block, MLP, and Attention, as shown in Table 1. Based on the observation that deeper layers are more redundant than shallow layers (Gromov et al., 2024; He et al., 2024b), we focus on deploying Router-Tuning to the deepest layers except the last one, leaving other layers unchanged.

While previous studies have primarily explored layer dropping or skipping to Block and MLP layers (Bae et al., 2023; Gromov et al., 2024), skipping these modules significantly degrades performance when applied at either token or sequence level. In

Table 1: **Router-Tuning at different granularities.** We compare deployments on Attention, Block, and MLP layers. The number of skippable layers is capped at 16, with 50% MoD capacity. SpeedUp denotes inference-time speedup.

Llama-3-8B											
Method	Granularity	Speedup	ARC-C	BoolQ	HellaSwag	MMLU	OBQA	PIQA	RTE	WinoGrande	Avg.
Baseline	–	1.00×	58.1	81.3	82.1	65.3	45.0	80.5	67.2	77.7	69.7
Router-Tuning	Block _{token}	1.24×	44.2	77.9	63.1	64.4	34.0	70.4	65.4	71.6	61.4
	Block _{seq}	1.26×	44.5	78.0	62.6	64.6	34.2	70.3	65.3	71.2	61.3
	MLP _{token}	1.05×	45.3	77.8	65.1	62.8	33.7	71.9	66.8	72.4	62.0
	MLP _{seq}	1.06×	45.1	77.7	65.4	62.4	33.4	71.6	66.4	72.1	61.8
	Attn _{token}	1.18×	56.4	79.8	81.0	65.3	45.2	79.9	64.6	77.3	68.7
	Attn _{seq}	1.21×	56.6	80.5	80.7	65.1	44.6	80.5	69.7	77.7	69.4
Llama-3-8B-Instruct											
Method	Granularity	Speedup	ARC-C	BoolQ	HellaSwag	MMLU	OBQA	PIQA	RTE	WinoGrande	Avg.
Baseline	–	1.00×	62.1	83.2	78.8	65.7	42.8	78.7	67.5	75.9	69.3
Router-Tuning	Block _{token}	1.24×	44.6	80.9	54.1	60.2	31.2	64.8	67.7	65.1	58.6
	Block _{seq}	1.26×	44.7	81.2	54.5	60.6	32.4	64.6	67.1	64.8	58.7
	MLP _{token}	1.05×	41.4	74.9	59.3	64.8	31.6	67.8	66.4	68.4	59.3
	MLP _{seq}	1.06×	41.8	75.1	59.3	64.5	31.2	68.2	66.7	68.8	59.5
	Attn _{token}	1.18×	60.2	82.9	76.8	65.8	42.6	78.6	67.7	76.6	68.9
	Attn _{seq}	1.21×	60.4	83.3	76.9	65.7	43.0	78.2	68.2	76.9	69.1

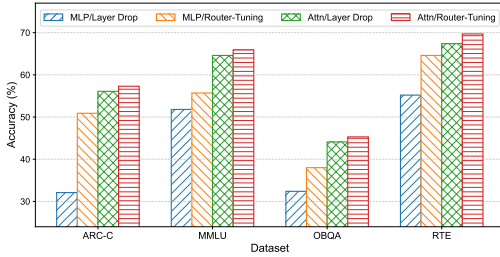


Figure 2: **Comparison between Router-Tuning and Layer Drop** on MLP and Attention layers under a fixed 25% overall skipping ratio.

contrast, applying dynamic depth to Attention layers maintains the performance of original models, e.g., 69.4 v.s. 69.7 in Llama-3-8B. These findings reinforce our motivation to target Attention layers, and we utilize Router-Tuning on Attention layers by default unless stated otherwise.

Router-Tuning improves over static layer dropping While statically dropping attention layers (He et al., 2024b) has demonstrated promising performance, its static nature lacks flexibility and limits its representational power. Here, we further investigate the improvements offered by the dynamic mechanism. Figure 2 compares Router-Tuning with static Layer Drop (He et al., 2024b), where Router-Tuning consistently achieves superior performance. For more complex tasks that are more sensitive to layer skipping, as shown in Figure 3, we compare Router-Tuning with Layer Drop on attention layers (i.e., “Attention Drop”) under the same computation budget, e.g., dropping 4 layers versus deploying MoD to 8 layers with 50% capacity. Under the same skipping ratios, Router-Tuning significantly outperforms Attention Drop on the GSM8K benchmark (Cobbe et al., 2021), e.g., 6.5% when the skipping ratio is 25.0%. In Figure 4, we further visualize the layer-wise skipping

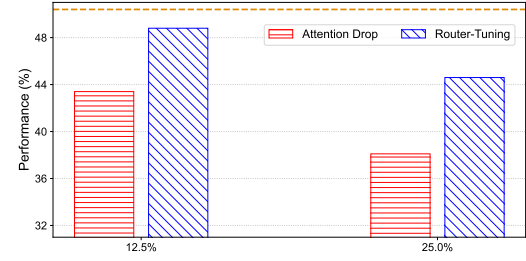


Figure 3: **Comparison with Attention Drop** on GSM8K tasks under identical skipping ratios.

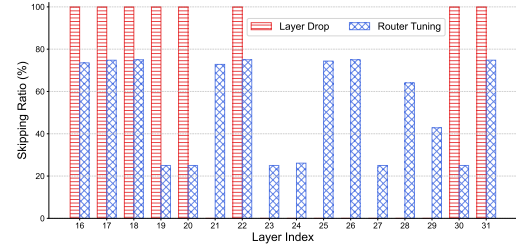


Figure 4: **Layer-wise Skipping Ratios** for Attention layers after Layer Drop and Router-Tuning.

ratios of MoD versus Attention Drop. Unlike static approaches that permanently remove certain layers, Router-Tuning maintains the utilization of all layers by adaptively distributing skipping ratios across them. This flexible allocation strategy contributes to improved performance.

Router-Tuning outperforms dynamic skipping methods Table 2 compares Router-Tuning with dynamic skipping baselines, including DLO (Tan et al., 2024) and Skip Transformer (Peroni and Bertsimas, 2024). Following the setup of baselines, we conduct the comparison on the token level for MLP layers. Router-Tuning consistently outperforms these methods across most tasks, despite operating under the router-only training constraint. On the one hand, Router-Tuning freezes the model backbones, which contributes to avoiding the risks

of catastrophic forgetting (Houlsby et al., 2019b; Liu et al., 2024; Qiao and Mahdavi, 2024). Additionally, Router-Tuning is trained end-to-end without being constrained by precomputed labels (e.g., token-level similarity scores in DLO) or stochastic gating mechanisms. This design enables more flexible and stable learning, while preserving the pretrained capabilities of the backbone.

Table 2: **Performance comparison against dynamic dropping baselines**, including DLO (Tan et al., 2024) and Skip Transformer (Peroni and Bertsimas, 2024).

Method	ARC-C	HellaSwag	MMLU	WinoG	Avg.
DLO	44.5	64.2	62.1	71.3	60.5
Skip Transformer	44.7	64.4	62.4	71.5	60.8
Router-Tuning	45.3	65.1	62.8	72.4	61.4

4.2 Efficiency Improvements

In this part, we measure the efficiency in both training and inference, focusing on computational and memory usage.

Table 3: **Comparison of training strategies** of achieving dynamic layer skipping. The training time for MoD is left blank as it was not conducted on LLaMA-3-8B.

Method	Target Modules	Granularity	Training Stage	Trainable	Training Time
MoD (Raposo et al., 2024)	Block	Token	Pretaining	Full Model	36h on NVIDIA A100
DLO (Tan et al., 2024)	MLP	Token	Continual Pretaining	Full Model	36h on NVIDIA A100
Router-Tuning	Block / MLP / Attn	Token / Sequence	Finetuning	Router	15m on NVIDIA A6000

Training Efficiency The training efficiency of our method lies in two perspectives: trainable parameters and training steps. Since the router projects the input from dimension d to 1, the number of trainable parameters is $d \times 1$ per layer, and the total number of trainable parameters is fewer than 0.01% of the whole model. Additionally, Router-Tuning only requires a few steps, which is verified in Section 5. Consequently, as shown in Table 3, Router-Tuning can be completed in under 15 minutes on a single NVIDIA A6000 GPU—over 1000 times faster than DLO (Tan et al., 2024), which performs large-scale training on full models and takes 36 hours on NVIDIA RTX A100 GPUs.

Inference Speedup We also evaluate the runtime speed improvements achieved with Router-Tuning. The inference speed is measured throughout the entire generation process, including prefilling and generation. To ensure that our results accurately reflect the performance gains, we adhere to two key principles: (1) all operations are performed on a single Nvidia RTX A6000 Ada GPU, eliminating any communication overhead from multi-GPU setups; and (2) we set the maximal sequence length

as 2048 and increase batch sizes to fully utilize the GPU for each model.

As shown in Table 1, skipping attention layers yields a more substantial speedup than skipping MLP layers, which is primarily due to the quadratic complexity of the attention mechanism and the memory overhead associated with KV-cache (Zhang et al., 2023; Singhania et al., 2024; He et al., 2024b). On the other hand, different granularities contribute to different levels of speedup. This variation is primarily due to attention layers, where fine-grained token-level MoD introduces differing token lengths within a batch, necessitating padding operations to standardize sequence lengths. Instead, the speedups at the sequence level surpass those at the token level, with Router-Tuning achieving a 21% improvement in inference speed. Therefore, we set Router-Tuning on the sequence level for attention layers as the default setting.

KV Cache The KV cache stores intermediate representations of attention layers, accelerating inference by eliminating redundant computations but incurring substantial memory overhead. Our approach, which selectively skips attention layers, significantly reduces KV cache size—for instance, achieving an 8GB reduction when processing an input sequence of length 2048 with a batch size of 64 on Llama-3-8B. In contrast, DLO (Tan et al., 2024) operates exclusively on MLP layers and retains the full KV cache, providing no memory savings.

5 Ablation Studies

Compatibility across different models Since Router-Tuning can be seamlessly integrated into pretrained language models, we extend our evaluation to diverse architectures, including Llama-2 (Touvron et al., 2023), Mistral (Jiang et al., 2023), and Qwen2.5 (Bai et al., 2023), covering a wide range of model sizes. As shown in Table 4, we deploy Router-Tuning in half of the attention layers while maintaining the total MoD capacity at 50%. Across all models, Router-Tuning preserves performance compared to their original counterparts, demonstrating its effectiveness and adaptability across different architectures.

Impact of number of MoD layers In the main experiments, we deployed half of the layers with MoD. In Figure 5, we further explore the effect of the number of MoD layers. Our results indicate that applying MoD to up to half of the attention layers

Table 4: **Ablation study of Router-Tuning across multiple model architectures and scales**, highlighting its robustness and consistent improvements.

Models	Speedup	OBQA	PIQA	RTE	WinoGrande	BoolQ	ARC-C	HellaSwag	MMLU	Avg.
Llama-2-13B	1.00×	45.2	80.5	65.0	76.2	80.7	59.4	82.2	54.6	<u>68.0</u>
w/Router-Tuning	1.22×	45.4	80.6	64.6	76.2	80.5	59.3	82.2	54.7	67.9
Qwen-2.5-14B	1.00×	45.6	82.2	79.1	80.4	85.3	67.2	84.3	79.7	<u>75.5</u>
w/Router-Tuning	1.18×	45.4	82.4	77.9	78.5	85.0	66.2	83.8	78.0	74.7
Qwen-2.5-7B	1.00×	47.2	79.6	81.2	76.6	84.6	63.7	80.2	74.1	<u>73.4</u>
w/Router-Tuning	1.19×	47.0	80.1	76.9	76.1	83.2	62.3	79.8	73.3	72.3
Mistral-7B	1.00×	44.4	82.2	68.2	79.0	82.2	60.6	83.2	62.4	<u>70.3</u>
w/Router-Tuning	1.24×	44.0	81.8	67.6	78.2	81.7	59.9	82.6	61.8	69.7

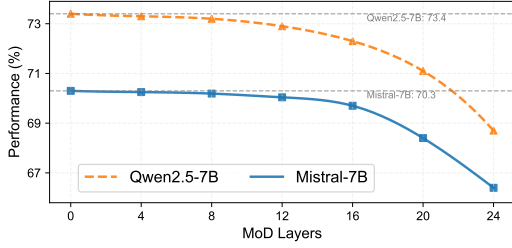


Figure 5: Ablation study on the impact of varying the number of MoD layers on overall model performance.

still maintains comparable performance. A similar trend is observed in Figure 6 for smaller models. However, when further increasing the number of MoD layers, performance starts to degrade. We attribute this decline to the transformation of important shallow layers, which negatively impacts overall performance (Men et al., 2024; He et al., 2024a,b). Therefore, preserving the density of shallow layers while applying MoD to deeper layers ensures its effectiveness.

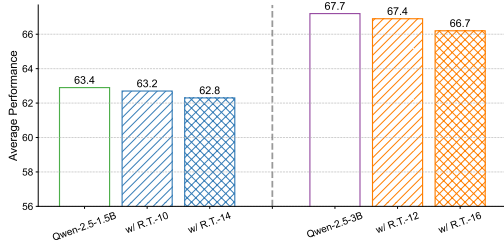


Figure 6: **Performance of Router-Tuning on small models**, where “R.T.” denotes Router-Tuning and the postfix “- n ” indicates that MoD is applied to n layers.

Influence of Training Dataset In Table 5, we next examine the impact of using different training datasets for Router-Tuning. We consider a variety of datasets, including Alpaca (Taori et al., 2023), Evol-Instruct (Xu et al., 2023), ShareGPT (Zheng et al., 2023), and Llama-Pro (Wu et al., 2024). Since Router-Tuning only fine-tunes the routers while keeping the backbone of the language models intact, changes in the training dataset do not significantly impact performance. However, Llama-Pro, which incorporates diverse training data from various domains, provides slightly better performance due to its broader data coverage.

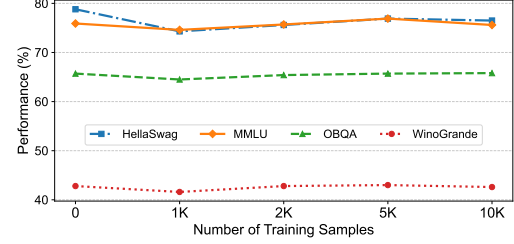


Figure 7: Effect of varying the number of training samples on performance.

On the other hand, due to the small number of trainable parameters, Router-Tuning does not require a large amount of training samples. As shown in Figure 7, MoD layers are initially dense-activated and then sparsified. Although the initial sparsification steps lead to a drop in performance, subsequent Router-Tuning facilitates performance recovery. Notably, just 5K training samples are sufficient to effectively train the routers.

Table 5: **Effectiveness across different training datasets**, where Router Tuning demonstrates robustness to the varying datasets.

Dataset	HellaSwag	MMLU	OBQA	WinoGrande	Avg.
Baseline	82.1	65.3	45.0	77.7	<u>67.5</u>
Alpaca	79.8	62.2	43.8	77.4	<u>65.8</u>
Evol-Instruct	80.4	64.0	44.4	77.6	<u>66.6</u>
ShareGPT	80.6	63.3	45.4	76.7	<u>66.5</u>
Llama-Pro	80.7	65.1	44.6	77.7	67.0

6 MoE and LoRA Integration

In this section, we further explore the integration of Router-Tuning with other architectures and training techniques. First, we implement Router-Tuning on mainstream MoE architectures. Then, we combine Router-Tuning with LoRA fine-tuning to enhance both efficiency and performance.

Router-Tuning on MoE The Expert’s redundancy in MoE has been widely demonstrated in recent works (Lu et al., 2024; He et al., 2024a), e.g., the model still maintains comparable performance after removing certain layers. Therefore, we further extend Router-Tuning to MoE, where we take OLMoE (Muennighoff et al., 2024) and DeepSeek-

Table 6: **Performance of Router-Tuning on Mixture of Experts**, where we take Expert Drop (He et al., 2024a) as the baseline of static dropping for comparison.

Models	SpeedUp	OBQA	PIQA	RTE	WinoGrande	BoolQ	ARC-C	HellaSwag	MMLU	Avg.
DeepSeek-MoE	1.00×	43.6	80.5	62.8	73.4	72.4	52.7	79.9	44.5	<u>63.7</u>
w/Expert Drop	1.11×	42.2	80.2	59.9	70.0	74.0	48.1	75.6	38.9	<u>61.1</u>
w/Router-Tuning	1.10×	43.2	80.4	61.2	71.4	72.1	50.8	77.3	42.8	62.4
OLMoE	1.00×	45.6	80.1	53.7	71.2	74.7	54.5	79.4	52.5	<u>64.0</u>
w/Expert Drop	1.13×	34.0	66.6	51.6	59.3	67.6	39.0	40.5	42.7	<u>50.2</u>
w/Router-Tuning	1.12×	40.4	76.2	53.2	70.2	71.3	52.0	77.9	50.1	61.4

Table 7: **Effectiveness of Router-Tuning integrated with LoRA finetuning**, compared to deploying Router-Tuning and LoRA separately.

Model	Method	SpeedUp	OBQA	PIQA	RTE	WinoGrande	BoolQ	ARC-C	HellaSwag	MMLU	Avg.
Llama-3-8B	Baseline	1.00×	45.0	80.5	67.2	77.7	81.3	58.1	82.1	65.3	<u>69.6</u>
	R.T.	1.21×	44.6	80.5	69.7	77.7	80.7	56.6	80.7	65.1	<u>69.5</u>
	LoRA	1.00×	46.6	82.0	68.0	77.9	83.9	61.8	81.6	65.9	71.0
	LoRA + R.T.	1.21×	47.2	82.2	67.4	77.8	83.9	61.5	81.7	65.8	<u>70.9</u>
Mistral-7B	Baseline	1.00×	44.4	82.2	68.2	79.1	82.2	60.6	83.2	62.4	<u>70.3</u>
	R.T.	1.24×	44.2	81.9	68.5	78.6	81.7	60.4	82.5	61.8	<u>70.0</u>
	LoRA	1.00×	45.2	83.0	68.9	79.4	84.7	60.7	83.7	62.8	71.1
	LoRA + R.T.	1.24×	45.7	83.1	68.7	79.3	84.3	60.9	83.4	62.9	71.2

MoE (Dai et al., 2024) as the backbones and equip each Expert network with Router-Tuning. Since these models deploy MoE at the token level, we directly apply the token-level Router-Tuning to these models. Here, we compare Router-Tuning with Expert Drop (He et al., 2024a) that statically drops less important experts. To ensure a fair comparison, we fix the overall skipping ratio of Router-Tuning to 25%, and drop the bottom 25% of experts globally by importance score in the Expert Drop baseline.

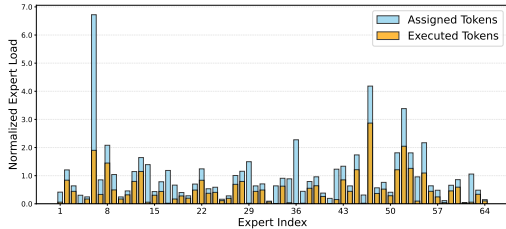


Figure 8: **Normalized expert load** before and after applying router-based filtering, denoted as “Assigned Tokens” and “Executed Tokens”, respectively. Expert load values are normalized by the mean number of assigned tokens.

In Table 6, instead of removing a subset of experts like Expert Drop, Router-Tuning maintains the potential of all experts, which contributes to a superior performance. In Figure 8, we further investigate how Router-Tuning affects the inference behavior of expert networks by visualizing the expert load within a specific layer from two perspectives: the number of tokens initially assigned to each expert (“Assigned Tokens”) and the number of tokens that pass the router and are actually executed (“Executed Tokens”). Router-tuning prompts the experts to skip less important tokens and significantly lower the load of overloaded experts, which alleviates the imbalanced distribution of token as-

signments. Consequently, this approach fosters a more balanced utilization across the entire set of experts (Fedus et al., 2022; He et al., 2025) and thus enhances the efficiency.

Integration with LoRA Router-Tuning enables dynamic depth to improve computational efficiency, whereas parameter-efficient fine-tuning (PEFT) methods aim to update a small subset of parameters to enhance downstream task performance. To examine whether Router-Tuning is complementary to PEFT, we propose jointly conducting Router-Tuning with LoRA fine-tuning (Hu et al., 2021), targeting improvements in both efficiency and task performance. As shown in Table 7, this joint training strategy preserves the efficiency benefits of Router-Tuning while maintaining the performance gains achieved by LoRA. Together, the integration of Router-Tuning and LoRA offers a more advanced fine-tuning paradigm that further enhances overall model capability.

7 Conclusion

In this work, we investigate the dynamic depth mechanism from both design and training perspectives. We propose Router-Tuning, which effectively implements dynamic depth by fine-tuning only a minimal number of parameters in just a few steps. Additionally, we explore Router-Tuning across a variety of modules and granularities to evaluate its effectiveness across a wide range of models and tasks. These advancements provide valuable insights and practical solutions for deploying dynamic depth and enhancing the efficiency of large language models.

Limitations

Despite the progress achieved in this work, several limitations remain. First, while we have advanced MoD through Router-Tuning, other, potentially more sophisticated training strategies may further improve performance and merit future investigation. Second, due to computational resource constraints, our experiments were limited to a small set of models and tasks. Extending this approach to a broader range of architectures and applications would provide deeper insight into its generalizability and full potential.

References

2019. Winogrande: An adversarial winograd schema challenge at scale.
- Sangmin Bae, Jongwoo Ko, Hwanjun Song, and Se-Young Yun. 2023. Fast and robust early-exiting framework for autoregressive language models with synchronized parallel decoding. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5910–5924.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. *Qwen technical report*. Preprint, arXiv:2309.16609.
- Yoshua Bengio, Nicholas Léonard, and Aaron Courville. 2013. *Estimating or propagating gradients through stochastic neurons for conditional computation*. Preprint, arXiv:1308.3432.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2019. *Piqa: Reasoning about physical commonsense in natural language*. Preprint, arXiv:1911.11641.
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. *Boolq: Exploring the surprising difficulty of natural yes/no questions*. Preprint, arXiv:1905.10044.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. *Think you have solved question answering? try arc, the ai2 reasoning challenge*. Preprint, arXiv:1803.05457.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- 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. 2024. *Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models*. CoRR, 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, Miaojuan 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, Wanbiao 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, Yudian 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,

696	Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi	Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar	757
697	Gao, and Zizheng Pan. 2024. Deepseek-v3 technical	Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew	758
698	report . <i>Preprint</i> , arXiv:2412.19437.	Oldham, Mathieu Rita, Maya Pavlova, Melanie Kam-	759
699	Mostafa Elhoushi, Akshat Shrivastava, Diana Liskovich,	badur, Mike Lewis, Min Si, Mitesh Kumar Singh,	760
700	Basil Hosmer, Bram Wasti, Liangzhen Lai, Anas	Mona Hassan, Naman Goyal, Narjes Torabi, Niko-	761
701	Mahmoud, Bilge Acun, Saurabh Agarwal, Ahmed	lay Bashlykov, Nikolay Bogoychev, Niladri Chatterji,	762
702	Roman, Ahmed A Aly, Beidi Chen, and Carole-Jean	Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick	763
703	Wu. 2024. Layerskip: Enabling early exit inference	Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasić,	764
704	and self-speculative decoding .	Peter Weng, Prajwal Bhargava, Pratik Dubal,	765
705	William Fedus, Barret Zoph, and Noam Shazeer. 2022.	Praveen Krishnan, Punit Singh Koura, Puxin Xu,	766
706	Switch transformers: Scaling to trillion parameter	Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj	767
707	models with simple and efficient sparsity . <i>Preprint</i> ,	Ganapathy, Ramon Calderer, Ricardo Silveira Cabral,	768
708	arXiv:2101.03961.	Robert Stojnic, Roberta Raileanu, Rohan Maheswari,	769
709	Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman,	Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ron-	770
710	Sid Black, Anthony DiPofi, Charles Foster, Laurence	nie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan	771
711	Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li,	Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sa-	772
712	Kyle McDonell, Niklas Muennighoff, Chris Ociepa,	hana Chennabasappa, Sanjay Singh, Sean Bell, Seo-	773
713	Jason Phang, Laria Reynolds, Hailey Schoelkopf,	hyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sha-	774
714	Aviya Skowron, Lintang Sutawika, Eric Tang, An-	ran Narang, Sharath Raparthy, Sheng Shen, Shengye	775
715	ish Thite, Ben Wang, Kevin Wang, and Andy Zou.	Wan, Shruti Bhosale, Shun Zhang, Simon Van-	776
716	2023. A framework for few-shot language model	denhede, Soumya Batra, Spencer Whitman, Sten	777
717	evaluation .	Sootla, Stephane Collot, Suchin Gururangan, Syd-	778
718	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	ney Borodinsky, Tamar Herman, Tara Fowler, Tarek	779
719	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	Sheasha, Thomas Georgiou, Thomas Scialom, Tobias	780
720	Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-	Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal	781
721	ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh	Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh	782
722	Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-	Ramanathan, Viktor Kerkez, Vincent Gonguet, Vir-	783
723	tra, Archie Sravankumar, Artem Korenev, Arthur	ginie Do, Vish Voleti, Vitor Albiero, Vladan Petro-	784
724	Hinsvark, Arun Rao, Aston Zhang, Aurelien Ro-	vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whit-	785
725	driguez, Austen Gregerson, Ava Spataru, Baptiste	ney Meers, Xavier Martinet, Xiaodong Wang, Xi-	786
726	Roziere, Bethany Biron, Binh Tang, Bobbie Chern,	aofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinf-	787
727	Charlotte Caucheteux, Chaya Nayak, Chloe Bi,	feng Xie, Xuchao Jia, Xuewei Wang, Yaelle Gold-	788
728	Chris Marra, Chris McConnell, Christian Keller,	schlag, Yashesh Gaur, Yasmine Babaei, Yi Wen,	789
729	Christophe Touret, Chunyang Wu, Corinne Wong,	Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao,	790
730	Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-	Zacharie Delpierre Coudert, Zheng Yan, Zhengxing	791
731	lonsius, Daniel Song, Danielle Pintz, Danny Livshits,	Chen, Zoe Papakipos, Aaditya Singh, Aayushi Sri-	792
732	Danny Wyatt, David Esiobu, Dhruv Choudhary,	vastava, Abha Jain, Adam Kelsey, Adam Shajnfeld,	793
733	Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,	Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,	794
734	Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy,	Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei	795
735	Elina Lobanova, Emily Dinan, Eric Michael Smith,	Baevski, Allie Feinstein, Amanda Kallet, Amit San-	796
736	Filip Radenovic, Francisco Guzmán, Frank Zhang,	gani, Amos Teo, Anam Yunus, Andrei Lupu, And-	797
737	Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis An-	res Alvarado, Andrew Caples, Andrew Gu, Andrew	798
738	derson, Govind Thattai, Graeme Nail, Gregoire Mi-	Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchan-	799
739	alon, Guan Pang, Guillem Cucurell, Hailey Nguyen,	dani, Annie Dong, Annie Franco, Anuj Goyal, Apar-	800
740	Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan	jita Saraf, Arkabandhu Chowdhury, Ashley Gabriel,	801
741	Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Is-	Ashwin Bharambe, Assaf Eisenman, Azadeh Yaz-	802
742	han Misra, Ivan Evtimov, Jack Zhang, Jade Copet,	dan, Beau James, Ben Maurer, Benjamin Leonhardi,	803
743	Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park,	Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi	804
744	Jay Mahadeokar, Jeet Shah, Jelmer van der Linde,	Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Han-	805
745	Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu,	cock, Bram Wasti, Brandon Spence, Brani Stojkovic,	806
746	Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang,	Brian Gamido, Britt Montalvo, Carl Parker, Carly	807
747	Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park,	Burton, Catalina Mejia, Ce Liu, Changhan Wang,	808
748	Joseph Rocca, Joshua Johnstun, Joshua Saxe, Jun-	Changkyu Kim, Chao Zhou, Chester Hu, Ching-	809
749	teng Jia, Kalyan Vasuden Alwala, Karthik Prasad,	Hsiang Chu, Chris Cai, Chris Tindal, Christoph Fe-	810
750	Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth	ichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty,	811
751	Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer,	Daniel Kreymer, Daniel Li, David Adkins, David	812
752	Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal	Xu, Davide Testuggine, Delia David, Devi Parikh,	813
753	Lakhotia, Lauren Rantala-Yearly, Laurens van der	Diana Liskovich, Didem Foss, Dingkan Wang, Duc	814
754	Maaten, Lawrence Chen, Liang Tan, Liz Jenkins,	Le, Dustin Holland, Edward Dowling, Eissa Jamil,	815
755	Louis Martin, Lovish Madaan, Lubo Malo, Lukas	Elaine Montgomery, Eleonora Presani, Emily Hahn,	816
756	Blecher, Lukas Landzaat, Luke de Oliveira, Madeline	Emily Wood, Eric-Tuan Le, Erik Brinkman, Este-	817
		ban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun,	818
		Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat	819
		Ozgenel, Francesco Caggioni, Frank Kanayet, Frank	820

821	Seide, Gabriela Medina Florez, Gabriella Schwarz,	Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia,	885
822	Gada Badeer, Georgia Swee, Gil Halpern, Grant	Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,	886
823	Herman, Grigory Sizov, Guangyi, Zhang, Guna	Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao,	887
824	Lakshminarayanan, Hakan Inan, Hamid Shojanaz-	Yundi Qian, Yunlu Li, Yuze He, Zach Rait, Zachary	888
825	eri, Han Zou, Hannah Wang, Hanwen Zha, Haroun	DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang,	889
826	Habeeb, Harrison Rudolph, Helen Suk, Henry As-	Zhiwei Zhao, and Zhiyu Ma. 2024. The llama 3 herd	890
827	pegren, Hunter Goldman, Hongyuan Zhan, Ibrahim	of models . <i>Preprint</i> , arXiv:2407.21783.	891
828	Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis,		
829	Irina-Elena Veliche, Itai Gat, Jake Weissman, James	Andrey Gromov, Kushal Tirumala, Hassan Shapourian,	892
830	Geboski, James Kohli, Janice Lam, Japhet Asher,	Paolo Gloriosi, and Daniel A. Roberts. 2024. The	893
831	Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jen-	unreasonable ineffectiveness of the deeper layers .	894
832	nifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy	<i>Preprint</i> , arXiv:2403.17887.	895
833	Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe		
834	Cummings, Jon Carvill, Jon Shepard, Jonathan Mc-	Qi Han, Zejia Fan, Qi Dai, Lei Sun, Ming-Ming Cheng,	896
835	Phie, Jonathan Torres, Josh Ginsburg, Junjie Wang,	Jiaying Liu, and Jingdong Wang. 2022. On the con-	897
836	Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khan-	nection between local attention and dynamic depth-	898
837	delwal, Katayoun Zand, Kathy Matosich, Kaushik	wise convolution . In <i>International Conference on</i>	899
838	Veeraraghavan, Kelly Michelena, Keqian Li, Ki-	<i>Learning Representations</i> .	900
839	ran Jagadeesh, Kun Huang, Kunal Chawla, Kyle		
840	Huang, Lailin Chen, Lakshya Garg, Lavender A,	Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui	901
841	Leandro Silva, Lee Bell, Lei Zhang, Liangpeng	Wang, and Yulin Wang. 2021. Dynamic neural net-	902
842	Guo, Licheng Yu, Liron Moshkovich, Luca Wehrst-	works: A survey. <i>IEEE Transactions on Pattern Anal-</i>	903
843	edt, Madian Khabsa, Manav Avalani, Manish Bhatt,	<i>ysis and Machine Intelligence</i> , 44(11):7436–7456.	904
844	Martynas Mankus, Matan Hasson, Matthew Lennie,		
845	Matthias Reso, Maxim Groshev, Maxim Naumov,	Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-	905
846	Maya Lathi, Meghan Keneally, Miao Liu, Michael L.	Kirkpatrick, and Graham Neubig. 2021. Towards a	906
847	Seltzer, Michal Valko, Michelle Restrepo, Mihir Pa-	unified view of parameter-efficient transfer learning.	907
848	tel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark,	<i>arXiv preprint arXiv:2110.04366</i> .	908
849	Mike Macey, Mike Wang, Miquel Jubert Hermoso,		
850	Mo Metanat, Mohammad Rastegari, Munish Bansal,	Shwai He, Weilin Cai, Jiayi Huang, and Ang Li.	909
851	Nandhini Santhanam, Natascha Parks, Natasha	2025. Capacity-aware inference: Mitigating the	910
852	White, Navyata Bawa, Nayan Singhal, Nick Egebo,	straggler effect in mixture of experts . <i>Preprint</i> ,	911
853	Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich	arXiv:2503.05066.	912
854	Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz,		
855	Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin	Shwai He, Daize Dong, Liang Ding, and Ang Li.	913
856	Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pe-	2024a. Demystifying the compression of mixture-	914
857	dro Rittner, Philip Bontrager, Pierre Roux, Piotr	of-experts through a unified framework . <i>Preprint</i> ,	915
858	Dollar, Polina Zvyagina, Prashant Ratanchandani,	arXiv:2406.02500.	916
859	Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel		
860	Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu	Shwai He, Guoheng Sun, Zheyu Shen, and Ang Li.	917
861	Nayani, Rahul Mitra, Rangaprabhu Parthasarathy,	2024b. What matters in transformers? not all atten-	918
862	Raymond Li, Rebekkah Hogan, Robin Battey, Rocky	tion is needed . <i>Preprint</i> , arXiv:2406.15786.	919
863	Wang, Russ Howes, Ruty Rinott, Sachin Mehta,		
864	Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	920
865	Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov,	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	921
866	Satadru Pan, Saurabh Mahajan, Saurabh Verma,	2021. Measuring massive multitask language under-	922
867	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-	standing . <i>Preprint</i> , arXiv:2009.03300.	923
868	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,		
869	Shengxin Cindy Zha, Shishir Patil, Shiva Shankar,	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,	924
870	Shuqiang Zhang, Shuqiang Zhang, Sinong Wang,	Bruna Morrone, Quentin de Laroussilhe, Andrea	925
871	Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala,	Ges mundo, Mona Attariyan, and Sylvain Gelly.	926
872	Stephanie Max, Stephen Chen, Steve Kehoe, Steve	2019a. Parameter-efficient transfer learning for nlp .	927
873	Satterfield, Sudarshan Govindaprasad, Sumit Gupta,	<i>Preprint</i> , arXiv:1902.00751.	928
874	Summer Deng, Sungmin Cho, Sunny Virk, Suraj		
875	Subramanian, Sy Choudhury, Sydney Goldman, Tal	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,	929
876	Remez, Tamar Glaser, Tamara Best, Thilo Koehler,	Bruna Morrone, Quentin de Laroussilhe, Andrea Ges-	930
877	Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim	mundo, Mona Attariyan, and Sylvain Gelly. 2019b.	931
878	Matthews, Timothy Chou, Tzook Shaked, Varun	Parameter-efficient transfer learning for nlp . <i>ArXiv</i> ,	932
879	Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai	abs/1902.00751.	933
880	Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad		
881	Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu,	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	934
882	Vladimir Ivanov, Wei Li, Wenchen Wang, Wen-	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	935
883	wen Jiang, Wes Bouaziz, Will Constable, Xiaocheng	Weizhu Chen. 2021. Lora: Low-rank adaptation of	936
884	Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo	large language models . <i>Preprint</i> , arXiv:2106.09685.	937

938	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan	Dettmers, Douwe Kiela, Ali Farhadi, Noah A. Smith,	995
939	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	Pang Wei Koh, Amanpreet Singh, and Hannaneh	996
940	Weizhu Chen. 2022. LoRA: Low-rank adaptation of	Hajishirzi. 2024. Olmoe: Open mixture-of-experts	997
941	large language models . In <i>International Conference</i>	language models . <i>Preprint</i> , arXiv:2409.02060.	998
942	<i>on Learning Representations</i> .		
943	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	999
944	sch, Chris Bamford, Devendra Singh Chaplot, Diego	Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	1000
945	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	man, Diogo Almeida, Janko Altschmidt, Sam Alt-	1001
946	laume Lample, Lucile Saulnier, L��lio Renard Lavaud,	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	1002
947	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-	1003
948	Thibaut Lavril, Thomas Wang, Timoth��e Lacroix,	ing Bao, Mohammad Bavarian, Jeff Belgum, Ir-	1004
949	and William El Sayed. 2023. Mistral 7b . <i>Preprint</i> ,	wan Bello, Jake Berdine, Gabriel Bernadett-Shapiro,	1005
950	arXiv:2310.06825.	Christopher Berner, Lenny Bogdonoff, Oleg Boiko,	1006
951	Albert Q. Jiang, Alexandre Sablayrolles, Antoine	Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-	1007
952	Roux, Arthur Mensch, Blanche Savary, Chris	man, Tim Brooks, Miles Brundage, Kevin Button,	1008
953	Bamford, Devendra Singh Chaplot, Diego de las	Trevor Cai, Rosie Campbell, Andrew Cann, Brittany	1009
954	Casas, Emma Bou Hanna, Florian Bressand, Gi-	Carey, Chelsea Carlson, Rory Carmichael, Brooke	1010
955	anna Lengyel, Guillaume Bour, Guillaume Lam-	Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully	1011
956	ple, L��lio Renard Lavaud, Lucile Saulnier, Marie-	Chen, Ruby Chen, Jason Chen, Mark Chen, Ben	1012
957	Anne Lachaux, Pierre Stock, Sandeep Subramanian,	Chess, Chester Cho, Casey Chu, Hyung Won Chung,	1013
958	Sophia Yang, Szymon Antoniak, Teven Le Scao,	Dave Cummings, Jeremiah Currier, Yunxing Dai,	1014
959	Th��ophile Gervet, Thibaut Lavril, Thomas Wang,	Cory Decareaux, Thomas Degry, Noah Deutsch,	1015
960	Timoth��e Lacroix, and William El Sayed. 2024. Mix-	Damien Deville, Arka Dhar, David Dohan, Steve	1016
961	tral of experts . <i>Preprint</i> , arXiv:2401.04088.	Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti,	1017
962	Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang,	Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,	1018
963	Wei-Ming Chen, Wei-Chen Wang, Guangxuan	Sim��n Posada Fishman, Juston Forte, Isabella Ful-	1019
964	Xiao, Xingyu Dang, Chuang Gan, and Song Han.	ford, Leo Gao, Elie Georges, Christian Gibson, Vik	1020
965	2024. Awq: Activation-aware weight quantization	Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-	1021
966	for llm compression and acceleration . <i>Preprint</i> ,	Lopes, Jonathan Gordon, Morgan Grafstein, Scott	1022
967	arXiv:2306.00978.	Gray, Ryan Greene, Joshua Gross, Shixiang Shane	1023
968	Shuo Liu, Jacky Keung, Zhen Yang, Fang Liu, Qilin	Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris,	1024
969	Zhou, and Yihan Liao. 2024. Delving into parameter-	Yuchen He, Mike Heaton, Johannes Heidecke, Chris	1025
970	efficient fine-tuning in code change learning: An	Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele,	1026
971	empirical study . <i>Preprint</i> , arXiv:2402.06247.	Brandon Houghton, Kenny Hsu, Shengli Hu, Xin	1027
972	Xudong Lu, Qi Liu, Yuhui Xu, Aojun Zhou, Siyuan	Hu, Joost Huizinga, Shantanu Jain, Shawn Jain,	1028
973	Huang, Bo Zhang, Junchi Yan, and Hongsheng Li.	Joanne Jang, Angela Jiang, Roger Jiang, Haozhun	1029
974	2024. Not all experts are equal: Efficient expert	Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee-	1030
975	pruning and skipping for mixture-of-experts large	woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-	1031
976	language models . In <i>Proceedings of the 62nd Annual</i>	mali, Ingmar Kanitscheider, Nitish Shirish Keskar,	1032
977	<i>Meeting of the Association for Computational Lin-</i>	Tabarak Khan, Logan Kilpatrick, Jong Wook Kim,	1033
978	<i>guistics (Volume 1: Long Papers)</i> , pages 6159–6172,	Christina Kim, Yongjik Kim, Jan Hendrik Kirchner,	1034
979	Bangkok, Thailand. Association for Computational	Jamie Kiros, Matt Knight, Daniel Kokotajlo,	1035
980	Linguistics.	Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-	1036
981	Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang,	stantinidis, Kyle Kopic, Gretchen Krueger, Vishal	1037
982	Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	1038
983	Chen. 2024. Shortgpt: Layers in large language mod-	Leike, Jade Leung, Daniel Levy, Chak Ming Li,	1039
984	els are more redundant than you expect . <i>Preprint</i> ,	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	1040
985	arXiv:2403.03853.	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,	1041
986	Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish	Anna Makanju, Kim Malfacini, Sam Manning, Todor	1042
987	Sabharwal. 2018. Can a suit of armor conduct elec-	Markov, Yaniv Markovski, Bianca Martin, Katie	1043
988	tricity? a new dataset for open book question answer-	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer	1044
989	ing . <i>Preprint</i> , arXiv:1809.02789.	McKinney, Christine McLeavey, Paul McMillan,	1045
990	Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld,	Jake McNeil, David Medina, Aalok Mehta, Jacob	1046
991	Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi,	Menick, Luke Metz, Andrey Mishchenko, Pamela	1047
992	Pete Walsh, Oyvind Tafjord, Nathan Lambert, Yuling	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	1048
993	Gu, Shane Arora, Akshita Bhagia, Dustin Schwenk,	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	1049
994	David Wadden, Alexander Wettig, Binyuan Hui, Tim	M��ly, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,	1050
		Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	1051
		Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	1052
		Paino, Joe Palermo, Ashley Pantuliano, Giambat-	1053
		tista Parascandolo, Joel Parish, Emy Parparita, Alex	1054
		Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	1055
		man, Filipe de Avila Belbute Peres, Michael Petrov,	1056
		Henrique Ponde de Oliveira Pinto, Michael, Poko-	1057

1058	rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report . <i>Preprint</i> , arXiv:2303.08774.	
1087	Matthew Peroni and Dimitris Bertsimas. 2024. Skip transformers: Efficient inference through skip-routing . In <i>NeurIPS 2024 Workshop on Fine-Tuning in Modern Machine Learning: Principles and Scalability</i> .	
1092	Fuli Qiao and Mehrdad Mahdavi. 2024. Learn more, but bother less: parameter efficient continual learning . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	
1096	David Raposo, Sam Ritter, Blake Richards, Timothy Lillicrap, Peter Conway Humphreys, and Adam Santoro. 2024. Mixture-of-depths: Dynamically allocating compute in transformer-based language models . <i>Preprint</i> , arXiv:2404.02258.	
1101	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer . <i>Preprint</i> , arXiv:1701.06538.	
1106	Prajwal Singhania, Siddharth Singh, Shwai He, Soheil Feizi, and Abhinav Bhatle. 2024. Loki: Low-rank keys for efficient sparse attention . <i>ArXiv</i> , abs/2406.02542.	
1110	Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. 2024. A simple and effective pruning approach for large language models . <i>Preprint</i> , arXiv:2306.11695.	
1113	Zhen Tan, Daize Dong, Xinyu Zhao, Jie Peng, Yu Cheng, and Tianlong Chen. 2024. Dlo: Dynamic layer operation for efficient vertical scaling of llms . <i>Preprint</i> , arXiv:2407.11030.	
1117	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca .	1119
1122	Gemini Team. 2024. Gemini 1.5: Unlocking multi-modal understanding across millions of tokens of context . <i>Preprint</i> , arXiv:2403.05530.	1123
1125	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>Preprint</i> , arXiv:2307.09288.	1124
1148	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In the Proceedings of ICLR.	1149
1153	Chengyue Wu, Yukang Gan, Yixiao Ge, Zeyu Lu, Jiahao Wang, Ye Feng, Ying Shan, and Ping Luo. 2024. Llama pro: Progressive llama with block expansion . <i>Preprint</i> , arXiv:2401.02415.	1154
1157	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions . <i>Preprint</i> , arXiv:2304.12244.	1158
1162	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? <i>Preprint</i> , arXiv:1905.07830.	1163
1166	Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuan-dong Tian, Christopher Re, Clark Barrett, Zhangyang Wang, and Beidi Chen. 2023. H2o: Heavy-hitter oracle for efficient generative inference of large language models . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1167
1173	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,	1174

1175 Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang,
1176 Joseph E. Gonzalez, and Ion Stoica. 2023. [Judg-](#)
1177 [ing llm-as-a-judge with mt-bench and chatbot arena.](#)
1178 *Preprint*, arXiv:2306.05685.

A Experiment Setup

Models We conduct experiments on Llama (Touvron et al., 2023; Grattafiori et al., 2024), Qwen (Bai et al., 2023), and Mistral (Jiang et al., 2023) due to their competitive performance and widespread adoption. Additionally, we leverage OLMoE (Muennighoff et al., 2024) and Deepseek-MoE (Dai et al., 2024) as the backbone to deploy Router-Tuning on the Mixture of Experts framework.

Datasets For the training dataset, we used Llama-Pro (Wu et al., 2024), given it spanning general instruction, math, and code for the SFT process and offering a wealth of instruction data with varying complexity levels. To evaluate model performance, we report normalized zero-shot or few-shot accuracy on the LM-Harness benchmark. The number of shots for each task is detailed in Table 8, which includes multiple tasks: ARC-C (Clark et al., 2018), BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2021), OBQA (Mihaylov et al., 2018), PIQA (Bisk et al., 2019), RTE (Wang et al., 2019), WinoGrande (ai2, 2019) and GSM8K (Cobbe et al., 2021). The evaluation code is based on EleutherAI’s LM Harness framework (Gao et al., 2023).

Table 8: **Experimental settings for evaluation tasks.** “Norm” refers to the normalization performed with respect to the length of the input.

Task	Number of few-shot	Metric
BoolQ	0	Accuracy
RTE	0	Accuracy
OBQA	0	Accuracy (Norm)
PIQA	0	Accuracy (Norm)
MMLU	5	Accuracy
WinoGrande	5	Accuracy
GSM8K	5	Exact Match
HellaSwag	10	Accuracy (Norm)
ARC-C	25	Accuracy (Norm)

Hyperparameters We set τ as 0.5, which corresponds to the midpoint of the sigmoid function. To ensure that training starts from dense models, we initialize \mathbf{W} to zero, ensuring that $\mathbf{R}(\mathbf{x}) \geq \tau$ initially, i.e., training from dense models. To achieve the desired MoD capacity, we perform a grid search over the learning rate from $\{1\text{e-}5, 2\text{e-}5, 5\text{e-}5, 1\text{e-}4, 2\text{e-}4\}$ and the scale factor λ from $\{0, 0.1, 0.01, 0.001\}$, respectively.