

From Dense to Dynamic: Token-Difficulty Driven MoEfication of Pre-Trained LLMs

Anonymous ACL submission

Abstract

001 Training large language models (LLMs) for
002 different inference constraints is computationally
003 expensive, limiting control over efficiency-
004 accuracy trade-offs. Moreover, once trained,
005 these models typically process tokens uniformly,
006 regardless of their complexity, leading to static
007 and inflexible behavior. In this paper, we
008 introduce a post-training optimization framework,
009 DynaMoE, that adapts a pre-trained dense LLM
010 to a token-difficulty-driven Mixture-of-Experts
011 model with minimal fine-tuning cost. This
012 adaptation makes the model dynamic, with
013 sensitivity control to customize the balance
014 between efficiency and accuracy. DynaMoE
015 features a token-difficulty-aware router that
016 predicts the difficulty of tokens and directs
017 them to the appropriate sub-networks or
018 experts, enabling larger experts to handle
019 more complex tokens and smaller experts to
020 process simpler ones. Our experiments
021 demonstrate that DynaMoE can generate a range
022 of adaptive model variants of the existing
023 trained LLM with a single fine-tuning step,
024 utilizing only $10B$ tokens, a minimal cost
025 compared to the base model’s training. Each
026 variant offers distinct trade-offs between
027 accuracy and performance. Compared to the
028 baseline post-training optimization framework,
029 Flextron, our method achieves similar
030 aggregated accuracy across downstream tasks,
031 despite using only $\frac{1}{9}$ th of their fine-tuning
032 cost.

1 Introduction

033 Large language models (LLMs) have significantly
034 advanced the field of natural language processing,
035 showcasing strong capabilities in addressing
036 complex tasks (Brown et al., 2020; Touvron et al.,
037 2023a; Wei et al., 2022). However, their large
038 size presents challenges, particularly in terms
039 of high memory and computational demands,
040 which can limit their deployment in resource-
041 constrained settings. To address this, LLMs
042 must be optimized for specific memory and
043 computational constraints

(Touvron et al., 2023b). However, designing
043 multi-billion-parameter models for every use
044 case is not cost-effective, as it demands
045 substantial training time, data, and resources.
046

047 Some prior works have focused on adapting
048 large LLMs for resource-constrained use cases
049 by distilling knowledge from larger models into
050 smaller ones (Hsieh et al., 2023) or pruning
051 model parameters to reduce computational
052 demands (Sun et al., 2024). While these
053 methods effectively enable the use of large
054 LLMs in low-resource scenarios, they often
055 lead to performance degradation and require
056 careful balancing between efficiency and
057 accuracy. Alternatively, other approaches
058 have investigated many-in-one LLM designs,
059 MatFormer (Devvrit et al., 2023) and
060 SortedNet (Valipour et al., 2024), to employ
061 multiple sub-networks within a single model
062 to accommodate different computational
063 budgets. These architectures use nested
064 structures integrated into the standard LLM
065 framework. However, they require non-
066 standard methodologies and significantly
067 longer, more resource-intensive training
068 processes, which can offset the intended
069 efficiency benefits.

067 Mixture-of-Experts (MoE) models (Shazeer
068 et al., 2017; Du et al., 2021; Fedus et al.,
069 2022; Zoph et al., 2022; He, 2024) have
070 emerged as a promising alternative to dense
071 models, offering improved efficiency by
072 sparsely activating select sub-modules or
073 experts. This selective activation enables
074 MoEs to achieve high performance while
075 using fewer computational resources during
076 inference. However, training MoEs from
077 scratch remains resource-intensive and each
078 expert becomes static, often requiring a
079 fixed compute budget irrespective of the
080 input complexity.

079 Flextron (Cai et al., 2024) explored a
080 post-training methodology by integrating the
081 MoE concept into a nested elastic structure
082 within the MLP layers, creating heterogeneous
083 experts of different sizes, selected by a
084 router conditioned on the

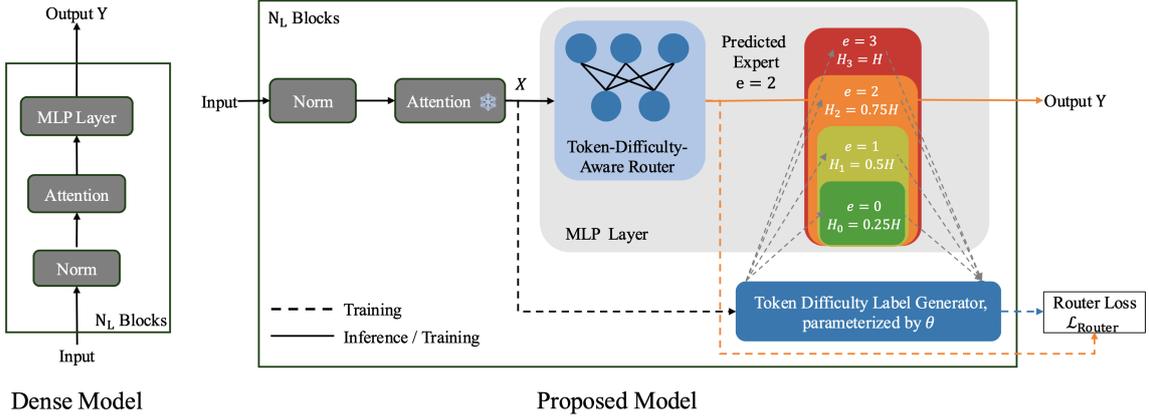


Figure 1: Overview of our proposed post-training optimization framework, DynaMoE. The left part represents the base pre-trained LLM, while the right part shows the adapted DynaMoE model.

input data. However, the lack of supervision in the router training leads to sub-optimal input complexity adaptation. Furthermore, the router lacks a parameter to customize its sensitivity to token complexity, limiting its flexibility and performance in handling diverse use-cases. Salehi et al. (2023) proposed an input-adaptive approach that predicts the difficulty of input data and dynamically adjusts the network’s width accordingly. In the absence of ground-truth difficulty labels, they relied on heuristic methods for label generation, which may limit precision and consistency in difficulty estimation.

To address their shortcomings, we introduce DynaMoE, a post-training optimization framework designed to transform a dense LLM into a token-difficulty-driven MoE model. DynaMoE leverages the insight that not all tokens require the full capacity of a model’s weights. For example, in the sentence “Geoffrey did his PhD at the university of Edinburgh”, simpler tokens like “at the university of” are predictable using prior context, while more complex tokens like “Edinburgh” demand broader contextual understanding. To maximize efficiency, DynaMoE selectively activates nested sub-components of the MLP, referred as experts, based on the predicted difficulty of each token. To this end, we make the following contributions:

- The framework includes a novel token-difficulty-aware router, trained to predict token hardness and assign it to the appropriate expert dynamically.
- Due to the lack of ground truth notion of hardness, we introduce a method to derive token difficulty labels which serve as supervision signals for training the router. This approach

allows a token to have varying difficulty labels across different layers.

- A simplified post-training optimization framework that efficiently adapts a pre-trained dense LLM into a token-difficulty-driven MoE model, featuring a sensitivity parameter to customize the efficiency vs accuracy trade-off.

2 Method

In this section, we describe our proposed post-training optimization framework, DynaMoE, which transforms a dense LLM into an MoE model for adaptive inference based on token difficulty. The process involves three key steps: (1) defining heterogeneous experts by splitting the MLP layers of the dense LLM; (2) generating token labels during training to represent token difficulty; and (3) training a router to predict token difficulty while fine-tuning the model. We detail these steps in the below sub-sections.

2.1 Defining Heterogeneous Experts

In this work, we focus on defining experts into the MLP layers of the LLM (Devvrit et al., 2023), as these layers account for the majority of the compute and operate on a token-by-token basis. The overview of DynaMoE is depicted in Fig. 1. The left part of the figure denotes the base pre-trained model which consists of the normalization layers, attention layers and the MLP layers in each transformer block. The right part shows the adapted DynaMoE model, where the original single MLP layer is transformed into multiple nested FFN blocks or experts. Such expert formation intro-

duces no additional parameters to the base model, aside from the router. This design draws inspiration from adaptive width reduction in transformer (Salehi et al., 2023) and recent works like Mat-former (Devvrit et al., 2023) and Flextron (Cai et al., 2024). The attention layers remain frozen, and the MLP layers adapted to nested experts are fine-tuned in DynaMoE.

Let D and H denote the embedding and the hidden dimensions of the MLP layer respectively. The input to the MLP layer is $\mathbf{X} \in \mathbb{R}^{B \times D}$ and the output is $\mathbf{Y} \in \mathbb{R}^{B \times D}$, where B is the batch dimension. The MLP layer with two fully connected layers is represented by weight matrices $\mathbf{W}^{(IN)} \in \mathbb{R}^{H \times D}$ and $\mathbf{W}^{(OUT)} \in \mathbb{R}^{D \times H}$. In order to get best results, we first rearrange these fully-connected layers, $\mathbf{W}^{(IN)}$ and $\mathbf{W}^{(OUT)}$, to have the most important rows/columns in the beginning of the matrix so that they can be included in all of the experts (Samragh et al., 2023). There are a total of E experts indexed using $e \in \{0, 1, \dots, E-1\}$. Each expert gets a portion H_e of the weight matrices $\mathbf{W}^{(IN)}$ and $\mathbf{W}^{(OUT)}$, sliced over the hidden dimension H . The value H_e is obtained as a fraction of H as,

$$H_e = \left\lfloor \left(\frac{e+1}{E} \right) \cdot H \right\rfloor, \quad (1)$$

consequently, $H_0 < H_1 < \dots < H_{E-1}$ and $H_{E-1} = H$. Note that the expert with index $E-1$ utilizes the full MLP layer. The restriction of the matrices $\mathbf{W}^{(IN)}$ and $\mathbf{W}^{(OUT)}$ to the expert width H_e is obtained using the slicing operator that selects the first H_e rows and columns respectively as

$$\mathbf{W}_e^{(IN)} = \mathbf{W}^{(IN)}[0 : H_e, :], \quad (2)$$

$$\mathbf{W}_e^{(OUT)} = \mathbf{W}^{(OUT)}[:, 0 : H_e]. \quad (3)$$

With σ as the activation function, the output \mathbf{Y}_e of the MLP layer corresponding to the expert e can thus be obtained as,

$$\mathbf{Y}_e = \sigma \left(\mathbf{X} \cdot \left(\mathbf{W}_e^{(IN)} \right)^T \right) \cdot \left(\mathbf{W}_e^{(OUT)} \right)^T. \quad (4)$$

2.2 Generating Token Difficulty Label

We aim to train a token-difficulty-aware router to dynamically assign tokens to an appropriate expert. But there is no ground-truth label denoting token difficulty to train such a router. To this end, we propose a method to estimate the token difficulty and generate a derived-ground-truth difficulty label

during training. This is shown as ‘‘Token Difficulty Label Generator’’ in Fig. 1.

First, we pass the input to all experts and generate the output \mathbf{Y}_e for each $e \in [E]$. Then, for each token $b \in [B]$ and each expert $e \in [E]$, we compute a similarity score $S_{b,e}$ that measures how similar is the output of the expert e compared to the output of the full MLP layer $e = E-1$ for that token. We calculate this similarity as,

$$S_{b,e} = \frac{\langle \mathbf{Y}_e[b, :], \mathbf{Y}_{E-1}[b, :] \rangle}{\langle \mathbf{Y}_{E-1}[b, :], \mathbf{Y}_{E-1}[b, :] \rangle}. \quad (5)$$

Here, $\langle \cdot, \cdot \rangle$ denotes the dot-product between two vectors. We use dot-product in this calculation as it accounts for both the magnitude and the direction of the tensors being compared.

Finally, we generate a derived ground-truth hardness label l_b , representing the target expert index for token b . Given a threshold θ , we assign l_b as the smallest expert index e satisfying $S_{b,e} > \theta$, that is, $l_b = \min\{e \in [E] \mid S_{b,e} > \theta\}$. We say that a token is easier if it has a smaller label l_b , that is the similarity score for a smaller expert is higher than threshold θ . In such cases, processing the token with the smaller expert incurs less compute without much compromise in the accuracy. During the forward pass of fine-tuning of DynaMoE, we generate the token difficulty labels. These labels are then used in the backward pass to compute the router loss, which trains the router.

2.3 Training a Token-Difficulty-Aware Router

The output of a router is in $\mathbb{R}^{B \times E}$, denoting logits over the E experts. Each router is parameterized by two linear layers, projecting the token embedding from dimension D to U and subsequently to E . In our experiments, we use $U = 256$, resulting in total parameters added by the routers across all layers to $33.6M$, which is only 0.51% of the base model size.

We train the router using the derived labels from Section 2.2. The objective is to learn the expert prediction using the derived-ground-truth labels to mimic token assignment based on their complexity and need. Hence, we impose the cross-entropy loss on the router to guide to this behavior and call it as the router loss. The overall objective function of DynaMoE is given as,

$$\mathcal{L} = \lambda_{LLM} \cdot \mathcal{L}_{LLM} + \lambda_{Router} \cdot \mathcal{L}_{Router}. \quad (6)$$

Here, \mathcal{L}_{LLM} is the main LLM Cross-entropy loss and \mathcal{L}_{Router} is the router loss. λ_{LLM} and λ_{Router}

	Cost (#Tokens)	Params	ARC-e	LAMBADA	PIQA	WinoGrande	Avg4	SciQ	HellaSwag	ARC-c	Avg7
Base Mistral 7B	-	7B	80.2	75.1	80.8	75.5	77.8	96.4	61.4	50.5	74.2
DynaMoE $\theta = 0.9$	10B	6B	75.0	71.0	78.3	71.8	74.0	95.2	57.9	41.5	70.1
DynaMoE $\theta = 0.8$	10B	5.1B	69.9	68.0	78.0	66.0	70.5	94.2	54.5	35.2	66.5
DynaMoE $\theta = 0.7$	10B	4.6B	66.0	65.9	75.4	63.4	67.7	93.6	52.1	31.7	64.0
Base Llama2-7B [†]	-	6.5B	75.1	71.5	77.5	69.1	73.3				
Flextron [†]	93.57B	4.1B	68.6	65.1	76.1	63.7	68.3				

Table 1: Evaluation of DynaMoE models with different sensitivity factor θ on downstream tasks, using 0-shot non-normalized accuracy metric. Our base model is Mistral 7B (Jiang et al., 2023). ([†]): results from Flextron (Cai et al., 2024) used as our baseline. *Params* denotes the average number of total activated parameters, aggregated over the downstream tasks. *Avg4* averages over *ARC-e*, *LAMBADA*, *PIQA*, *WinoGrande*, while *Avg7* averages over all tasks.

are hyper-parameters, denoting the weights of the respective losses.

3 Experimental Set up

3.1 Model and Dataset

Model: DynaMoE provides a simplified post-training approach to convert any dense LLM to an MoE model with a tunable sensitivity factor θ , the similarity threshold, to achieve desired latency reduction and the tolerance for drop in accuracy. DynaMoE integrates seamlessly with any transformer model, regardless of the architecture. To showcase the effectiveness of the method, we use Mistral 7B model (Jiang et al., 2023), a widely-used open-source pre-trained language model, as the base model.

Dataset: For DynaMoE fine-tuning, we use a small subset (10B tokens) of the Falcon Refined-Web dataset (Penedo et al., 2023). Falcon Refined-Web is an open-source dataset which contains high-quality web data. This minimal fine-tuning overhead enables a cost-effective conversion of any pre-trained LLM into an MoE variant for faster inference.

3.2 Training Details

We first reorder the pre-trained weight matrices (Samragh et al., 2023) in the MLP blocks before fine-tuning DynaMoE so that the most important weights can be included in all the experts. We run the dense LLM on a subset of the data and collect the absolute activations from the MLP layer, $Y \in R^{B,H}$. This subset can be a very small portion of the training data, 0.004% tokens of the Falcon RefinedWeb in our case. Subsequently, we aggregate the activations along the batch dimension and

across all the subset samples to obtain an importance score for each neuron in the hidden dimension. Using this importance score, we sort the MLP matrices.

Next, we fine-tune the DynaMoE model using only 10B tokens with AdamW optimizer (Loshchilov and Hutter, 2017) and a fixed learning rate of 10^{-5} . We keep the attention layers frozen. We set λ_{LLM} to 0.2 and λ_{Router} to 1 in Equation (6), the objective function for fine-tuning. We experiment with different values of threshold $\theta \in \{0.7, 0.8, 0.9\}$ to build a family of DynaMoE models with varying sensitivity parameter. A low sensitivity parameter, that is a smaller value of θ , makes the system less reactive, favoring smaller experts for most tokens and only escalating to bigger experts for significantly complex tokens. And a high sensitivity parameter makes the system more reactive, escalating to bigger experts even for moderately complex tokens. We use 4 experts ($E = 4$) with sizes $0.25H$, $0.5H$, $0.75H$, and H respectively. We denote the size of expert with index e as H_e .

4 Results

4.1 Evaluation

We evaluate the DynaMoE models on 7 downstream tasks using LM Evaluation Harness (Gao et al., 2024) and report the 0-shot non-normalized accuracy metric in Table 1. The selected evaluation tasks include ARC (Easy and Challenge) (Clark et al., 2018), HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2019), SciQ (Welbl et al., 2017), WinoGrande (Sakaguchi et al., 2019), and LAMBADA (Paperno et al., 2016).

DynaMoE is compared to two baselines, the

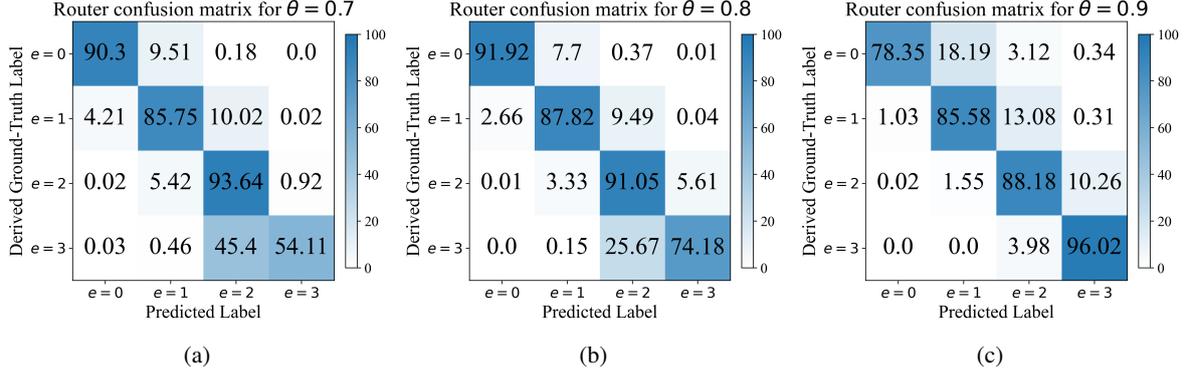


Figure 2: Confusion matrix for the router’s classification task in DynaMoE. The strong diagonal pattern in the matrices reflects high classification accuracy. The misclassified tokens often occurred in neighboring expert classes.

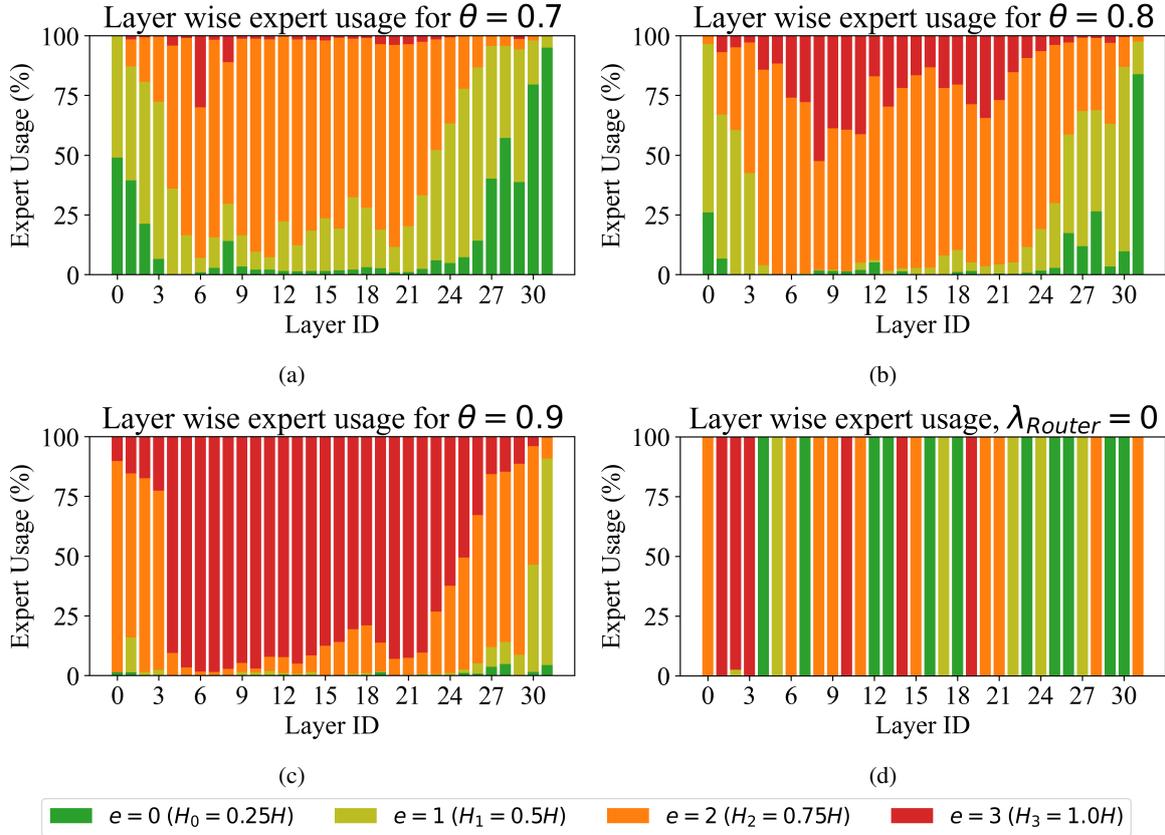


Figure 3: Layer-wise expert usage pattern of DynaMoE models with varying θ , aggregated across 7 downstream tasks. Each layer uses experts of varying sizes depending on the inputs. The sensitivity parameter θ regulates how readily tokens are routed to larger experts based on their difficulty: a lower θ favors smaller experts, while a higher θ prioritizes larger experts. 3d shows the case when DynaMoE is trained without the router loss λ_{Router} . In this case, the model tends to assign a specific expert per layer instead of dynamically selecting experts based on token difficulty.

base Mistral 7B model and the Flextron model (Cai et al., 2024) using Avg4 and Avg7 in Table 1. Compared to Mistral 7B, DynaMoE with $\theta = 0.8$ improves efficiency by activating only 5.1B of 7B parameters on average, with an 7.3 point accuracy drop after fine-tuning on only 10B tokens on the downstream tasks. The number of activated pa-

rameters adapts dynamically to token difficulty. For reference, Flextron fine-tunes on 93.57B tokens, activating 4.1B of 6.5B parameters, with a 5 point accuracy drop from its base model, Llama2-7B (Touvron et al., 2023a). We emphasize that with only $\frac{1}{9}$ th of the Flextron’s fine-tuning cost, our results for DynaMoE with $\theta = 0.7$ are com-

328	parable to Flextron. Accuracy improves with in-	using the router’s predicted expert indices without	378
329	crease in fine-tuning cost, but to keep the adaption	explicit difficulty supervision. The resulting expert	379
330	lightweight, we opt for a smaller cost.	usage pattern, shown in Fig. 3d, reveals that the	380
331	4.2 Analysis of Token-Difficulty-Aware	model converges to using specific experts per layer	381
332	Router	instead of dynamically allocating experts based on	382
333	We assess the performance of the Token-Difficulty-	token difficulty. In contrast, when router loss is	383
334	Aware router by gathering its predictions from all	applied, expert usage adapts dynamically to token	384
335	layers across the 7 downstream tasks outlined in	difficulty across layers.	385
336	Section 4.1. These predictions are then compared	5 Conclusion	386
337	to the ground truth labels derived in Section 2.2.	We present DynaMoE, a post-training optimiza-	387
338	Using both sets of labels, we compute the router’s	tion framework that converts a standard pre-trained	388
339	overall classification accuracy.	dense LLM into a token-difficulty-driven MoE	389
340	We present the confusion matrices for the	model. DynaMoE incorporates a lightweight	390
341	router’s classification tasks across all DynaMoE	router to predict the token difficulty and routes	391
342	models in Fig. 2. Notably, the matrices exhibit	them to an appropriate expert. To train this router,	392
343	a strong diagonal pattern, indicating high classi-	we propose a novel method to derive the token	393
344	fication accuracy. Furthermore, when the router	difficulty labels, which act as supervision signals.	394
345	misclassifies tokens, the errors predominantly oc-	DynaMoE generates adaptive model variants with	395
346	curred in neighboring expert classes, underscoring	sensitivity control, allowing customization of the	396
347	the router’s effectiveness in distinguishing token	trade-off between efficiency and accuracy.	397
348	difficulty levels.	Limitations	398
349	4.3 Experts usage analysis	While the proposed post-training optimization	399
350	We visualize the expert usage patterns across all	framework, DynaMoE, and the token-difficulty-	400
351	layers in Fig. 3. For each model, the Y-axis rep-	aware router provide a general-purpose approach	401
352	resents the percentage of tokens routed to a spe-	for token-difficulty adaptive processing, this work	402
353	cific expert, while the X-axis indicates the layer	does not explore their application to pre-trained	403
354	index. Notably, a token’s perceived difficulty may	heterogeneous mixture of expert (HMoE) models	404
355	vary across layers, hence it can be routed to dif-	(Wang et al., 2024). Such models, with their built-	405
356	ferent experts as the token progresses through the	in experts of varying capabilities, could benefit	406
357	model. The visualization shows that all experts	from the integration of DynaMoE, which provides	407
358	are utilized in varying proportions across layers,	a direct mechanism for routing tokens to appro-	408
359	reflecting an aggregated behavior over the 7 tasks.	prate experts based on token difficulty. Incorpor-	409
360	However, during inference, the model adapts to the	ating DynaMoE with HMoE models offers a	410
361	data, with simpler queries predominantly engaging	promising direction for future exploration. Addi-	411
362	lower-compute experts to maximize efficiency.	tionally, while DynaMoE makes the base model	412
363	The parameter θ affects expert usage in	input-adaptive, offering inference efficiency, this	413
364	DynaMoE models by controlling how quickly to-	work evaluates efficiency only in terms of the num-	414
365	kens are routed to larger experts based on diffi-	ber of active parameters and does not include the	415
366	culty. At lower θ values (e.g., $\theta = 0.7$), smaller	efficiency measured in a real deployment scenario.	416
367	experts ($e = 2$) dominate across layers, optimiz-	References	417
368	ing for efficiency. In contrast, at higher θ values	Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng	418
369	(e.g., $\theta = 0.9$), larger experts ($e = 3$) are uti-	Gao, and Yejin Choi. 2019. Piqa: Reasoning about	419
370	lized more frequently, prioritizing accuracy over	physical commonsense in natural language . In <i>AAAI</i>	420
371	efficiency. This shift demonstrates θ ’s role in bal-	<i>Conference on Artificial Intelligence</i> .	421
372	ancing computational resource allocation and pre-	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	422
373	diction accuracy.	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	423
374	Ablating the Router Loss: To examine the	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	424
375	router loss’s role in expert allocation, we train a	Askell, Sandhini Agarwal, Ariel Herbert-Voss,	425
376	DynaMoE model without it, relying solely on the		
377	LLM loss to train the router. Tokens are routed		

538 Hugo Touvron, Louis Martin, Kevin R. Stone, Peter
539 Albert, Amjad Almahairi, Yasmine Babaei, Niko-
540 lay Bashlykov, Soumya Batra, Prajjwal Bhargava,
541 Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cris-
542 tian Cantón Ferrer, Moya Chen, Guillem Cucurull,
543 David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin
544 Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami,
545 Naman Goyal, Anthony S. Hartshorn, Saghar Hos-
546 seini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor
547 Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V.
548 Korenev, Punit Singh Koura, Marie-Anne Lachaux,
549 Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai
550 Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov,
551 Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew
552 Poulton, Jeremy Reizenstein, Rashii Rungta, Kalyan
553 Saladi, Alan Schelten, Ruan Silva, Eric Michael
554 Smith, R. Subramanian, Xia Tan, Binh Tang, Ross
555 Taylor, Adina Williams, Jian Xiang Kuan, Puxin
556 Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, An-
557 gela Fan, Melanie Kambadur, Sharan Narang, Aure-
558 lien Rodriguez, Robert Stojnic, Sergey Edunov, and
559 Thomas Scialom. 2023b. [Llama 2: Open foundation
560 and fine-tuned chat models](#). *ArXiv*, abs/2307.09288.

561 Mojtaba Valipour, Mehdi Rezagholizadeh, Hossein Ra-
562 jabzadeh, Parsa Kavehzadeh, Marzieh Tahaei, Box-
563 ing Chen, and Ali Ghodsi. 2024. [Sortednet: A scal-
564 able and generalized framework for training modular
565 deep neural networks](#). *Preprint*, arXiv:2309.00255.

566 An Wang, Xingwu Sun, Ruobing Xie, Shuaipeng Li,
567 Jiaqi Zhu, Zhen Yang, Pinxue Zhao, J. N. Han, Zhan-
568 hui Kang, Di Wang, Naoaki Okazaki, Cheng zhong
569 Xu, and Tencent Hunyuan. 2024. [Hmoe: Hetero-
570 geneous mixture of experts for language modeling](#).
571 *ArXiv*, abs/2408.10681.

572 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
573 Bosma, Ed H. Chi, F. Xia, Quoc Le, and Denny Zhou.
574 2022. [Chain of thought prompting elicits reasoning
575 in large language models](#). *ArXiv*, abs/2201.11903.

576 Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017.
577 [Crowdsourcing multiple choice science questions](#).
578 *ArXiv*, abs/1707.06209.

579 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali
580 Farhadi, and Yejin Choi. 2019. [Hellaswag: Can a ma-
581 chine really finish your sentence?](#) In *Annual Meeting
582 of the Association for Computational Linguistics*.

583 Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du,
584 Yanping Huang, Jeff Dean, Noam Shazeer, and
585 William Fedus. 2022. [ST-MoE: designing stable and
586 transferable sparse expert models](#). *arXiv preprint
587 arXiv:2202.08906*.