Harder Task Needs More Experts: Dynamic Routing in MoE Models

Anonymous ACL submission

Abstract

In this paper, we introduce a novel dynamic expert selection framework for Mixture of Experts (MoE) models, aiming to enhance computational efficiency and model performance by adjusting the number of activated experts based on input difficulty. Unlike traditional MoE approaches that rely on fixed Top-K routing, which activates a predetermined number of experts regardless of the input's complexity, our method dynamically selects experts based on the confidence level in expert selection for 011 each input. This allows for a more efficient uti-012 lization of computational resources, activating 014 more experts for complex tasks requiring advanced reasoning and fewer for simpler tasks. Through extensive evaluations, our dynamic routing method demonstrates substantial improvements over conventional Top-2 routing 019 across various benchmarks, achieving an average improvement of 0.7% with less than 90% activated parameters. Further analysis shows our model dispatches more experts to tasks requiring complex reasoning skills, like BBH, confirming its ability to dynamically allocate computational resources in alignment with the input's complexity. Our findings also highlight a variation in the number of experts needed across different layers of the transformer model, offering insights into the potential for designing heterogeneous MoE frameworks. We will open-source all the models we trained in this project.

Introduction 1

034

To effectively increase the model's parameter size, researchers have proposed the Mixture of Experts (MoE) framework (Shazeer et al., 2017; Lepikhin et al., 2021). By setting up multiple experts to enhance the model's overall capacity, MoE models selectively activate a subset of parameters for use, thereby achieving more efficient parameter utilization. With the same number of activated parameters, MoE models substantially outperform 042

dense models in performance, achieving exceptional results in tasks such as QA and machine translation (Kim et al., 2021).

043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

Most MoE frameworks adopt a routing mechanism that dispatches a fixed number of experts for every input (Fedus et al., 2022; Du et al., 2022). The most famous method is Top-K routing (Shazeer et al., 2017), which initially calculates the probability of each expert being suited to the current input and then activates the Top-K suitable experts. Empirically, previous works (Lepikhin et al., 2021) activate two experts per token, as activating more experts offers limited improvements in model performance but substaintially increases training overhead. Most of the subsequent studies (Zoph et al., 2022; Lewis et al., 2021) can be seen as variants of Top-K routing, where different constraints are introduced to ensure that the number of tokens processed by different experts is as balanced as possible. Almost all these efforts activate a fixed number of experts.

The Top-K routing, though achieves good performance on downstream tasks, overlooks the different difficulties of inputs. Compared with simpler input, the more challenging input, e.g, tasks that require complex reasoning or logic inference, might need more parameters to solve. Dispatching experts equally across inputs could lead to computational waste on simpler tasks and insufficient computational resources for more difficult ones.

To fully leverage the potential of MoE models, we propose a dynamic routing mechanism that adjusts the number of required experts based on the confidence level in the expert selection. When the model deems the currently selected experts as insufficient, it activates more experts. Specifically, we first compute a probability distribution for selecting experts. If the highest probability for an expert exceeds a predefined threshold p, indicating high confidence, we activate only that one expert. Otherwise, we progressively include additional experts



Figure 1: Comparison between Top-K routing mechanism and Top-P routing mechanism. (a) Each token selects fixed K=2 experts with Top-K routing probabilities. (b) In Top-P routing mechanism, each token selects experts with higher routing probabilities until the cumulative probability exceeds threshold.

until the cumulative probability of the selected experts exceeds the threshold p. This approach allows for a dynamic selection of experts, with the number of experts adjusted according to the input's complexity.

Evaluation across multiple common benchmarks has revealed that our method substantially outperforms MoE models based on Top-K routing. Compared with Top-2 routing, our dynamic routing achieves an average improvement of 0.7% with less than 90% activated parameters. Further analysis has shown that our dynamic routing mechanism activates more experts in tasks requiring complex reasoning like BBH (Suzgun et al., 2023), while using fewer experts in relatively easier tasks such as Hellaswag (Zellers et al., 2019), confirming that our method indeed dynamically allocates experts based on the difficulty of the input. Token-level analysis indicates that tokens with ambiguous semantics are more challenging for the model, typically activating more experts. Another interesting finding is that the number of experts needed varies across different layers of the transformer. Lower layers require more experts for combination, while the top layer needs only one. This may relate to the *over-thinking* phenomenon (Kaya et al., 2019) widely observed in deep neural networks.

Our contributions can be summarized as follows:

- 1. We proposed a dynamic routing strategy that can adjust the number of activated experts based on the input difficulty dynamically.
- 2. We empirically validate that our proposed method is efficient in both training and inference, outperforming Top-2 routing while activating fewer experts.
- 3. We observe that for MoE models, the number

of experts needed to be activated varies across different layers. This finding could help design heterogeneous MoE frameworks.

2 Method

In this section, we first briefly introduce the MoE model with Top-K routing strategy, which activates a fixed number of experts for each token. As Top-K routing ignores the varying difficulty of different inputs and the different requirements for experts at different layers, we propose a dynamic routing mechanism that adjusts the number of activated experts according to the complexity of inputs. To avoid activating too many parameters through the dynamic routing mechanism, we also introduce a dynamic loss to encourage the model to activate only the necessary experts.

2.1 Top-K Routing MoE

In a Transformer model, the MoE layer is applied independently per token and replaces the feed-forward (FFN) sub-block of the transformer block (Lepikhin et al., 2021). For an MoE layer with N experts, $E = \{e_1, e_2, ..., e_N\}$, an input x will be sent to the experts and the output of the MoE layer is the weighted average of the experts':

$$MoE(\mathbf{x}) = \sum_{i=1}^{N} g_i(\mathbf{x}) * e_i(\mathbf{x})$$
(1)

where $g_*(\mathbf{x})$ is computed by a routing network that determines the contribution of each expert to the final output. In consideration of computing efficiency, a token is dispatched to limited experts. Thus for most experts, the corresponding $g_*(\mathbf{x})$ is zero meaning that the token is not dispatched to that expert.

107

108

109

110

111

112

113

114

115

116

118

123

120

121

122

124 125

126 127

128 129

130

131 132

133 134

135 136

138

139

140

141

142

143

144

145

146

147

148

149

151

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

198

199

To obtain $g_*(\mathbf{x})$, we first compute the probability **P** of selecting each expert for input **x**:

152

153

154

157

158

159

160

161

162

163

165

166

169

170

171

173

174

175

176

177

178

179

180

181

182

185

186

189

193

194

195

197

$$\mathbf{P} = Softmax(\mathbf{W}_{\mathbf{r}} \cdot \mathbf{x}^{T})$$
(2)

where $\mathbf{W}_{\mathbf{r}} \in N \times d$ is a learnable parameter and dis the dimension of the input \mathbf{x} . \mathbf{P} is a vector of size N and P_i represents the probability of selecting the i^{th} expert e_i to calculate the input \mathbf{x} .

Top-K routing selects the k experts, whose probabilities are the highest k in **P**. Then the probabilities of the selected experts are normalized and the weights of the remaining experts are set to zero, indicating they are not activated. The corresponding calculation of $g_*(\mathbf{x})$ is as follows:

$$g_i(\mathbf{x}) = \begin{cases} \frac{P_i}{\sum_{j \in TopK(\mathbf{P})} P_j}, & i \in TopK(\mathbf{P})\\ 0, & i \notin TopK(\mathbf{P}) \end{cases}$$
(3)

where $TopK(\mathbf{P})$ returns the indices of the largest k elements in \mathbf{P} .

Top-K routing is initially proposed by (Shazeer et al., 2017), and subsequently, numerous studies have built upon it with improvements. The following works (Lepikhin et al., 2021; Zuo et al., 2022) introduce constraints aimed at ensuring a more balanced workload among the experts during training. The core of these works remains to select the most suitable experts for each token under specific constraints, based on the probability distribution **P** calculated in Equation 2. And the number of experts dispatched for each token is fixed across all these studies. Empirically, the value of k is set to 2, serving as a trade-off between training costs and model capabilities.

2.2 Dynamic Routing MoE

Although the Top-K routing strategy has shown promising performance, its assumption that an equal number of experts should be dispatched for each token overlooks the variability in difficulty across different inputs. Moreover, since a fixed number of experts are activated at every layer of the transformer, this approach neglects the differences in representations across layers, potentially requiring a different number of experts for different layers.

To address these issues and make use of model parameters more efficiently, we propose a dynamic routing strategy based on model confidence. Unlike the Top-K routing, which selects a fixed number of experts, our method allows the model to assess whether the currently selected experts are sufficient. If not, it continues to incorporate more experts.

Specifically, we regard that P in Equation 2 reflects the confidence level of input \mathbf{x} in selecting different experts. In other words, P_i represents how confident the model is that the i^{th} expert can adequately handle input \mathbf{x} . If the highest probability in **P** is sufficiently large, then we may only need to use the corresponding expert. However, if the highest probability is not large enough, we need to add more experts to increase the reliability of processing x. We keep adding experts until the sum of the probabilities of the selected experts exceeds a specific threshold p, at which point we consider the model confident enough that these experts can effectively process the input x. We add new experts in descending order of their probabilities in P to minimize the number of activated experts as much as possible.

Formally, we first sort the elements in \mathbf{P} from highest to lowest, resulting in a sorted index list *I*. Then we find the smallest set of experts whose cumulative probability exceeds the threshold *p*, and the number of selected experts *t* is calculated by:

$$t = \underset{k \in \{1...,N\}}{\operatorname{argmin}} \sum_{j < =k} P_{i,j} \ge p \tag{4}$$

where p is the threshold that controls how confident the model should be when stopping adding more experts. p is a hyper-parameter whose range is from 0 to 1. The higher the p is, the more experts will be activated.

In dynamic routing mechanism, the calculation of $g_*(\mathbf{x})$ is:

$$g_i(\mathbf{x}) = \begin{cases} P_i & e_i \in S\\ 0, & e_i \notin S \end{cases}$$
(5)

where S is the set of selected experts controlled by t in Equation 4:

$$S = \{e_{I_1}, e_{I_2} \dots e_{I_t}\}$$
(6)

2.3 Loss

Dynamic Loss There is a risk associated with our dynamic routing mechanism: it could assign low confidence to all experts, thereby activating a larger number of experts to achieve better performance. Suppose \mathbf{P} is a uniform distribution and we set the hyper-parameter p to 0.5, then the model would activate up to half of the experts. This goes against

324

325

283

284

285

288

the original intention of the MoE framework, which is to scale the model with great efficiency.

To prevent dynamic routing from using too many parameters to cheat and losing its ability to selectively choose experts, we introduce a constraint on P. We expect the routing mechanism to select a small set of necessary experts, therefore, we aim to minimize the entropy of the distribution P, ensuring that every token can focus on as less specific experts as possible. Our dynamic loss is designed to encourage the routing mechanism to select the minimal necessary set of experts, which is formalized as:

$$Loss_d = -\sum_{i=1}^{N} P_i * log(P_i) \tag{7}$$

Load Balance Loss MoE models typically require distributed training, where different experts are deployed across various nodes. To avoid scenarios where some nodes are fully utilized while others are underutilized, thereby impacting training efficiency, it is generally desirable for the number of tokens processed by different experts to be roughly the same. Furthermore, previous study (Zuo et al., 2022) has shown that evenly activated experts in an MoE layer can lead to better performance. To achieve balanced loading among different experts, we have also incorporated a loadbalance loss, $Loss_b$, which is widely used in previous works (Lepikhin et al., 2021; Fedus et al., 2022)

$$Loss_b = N * \sum_{i=1}^{N} f_i * Q_i \tag{8}$$

where f_i is the fraction of the tokens choosing ex-272 pert e_i and Q_i is the fraction of the router probability allocated for expert e_i . For a sequence containing M tokens, f_i and Q_i are calculated as:

$$f_i = \frac{1}{M} \sum_{j=1}^M 1\{e_i \in S^j\}$$
(9)

277

278

242

243

244

245

246

247

251

256

257

260

261

262

266

271

273

$$Q_{i} = \frac{1}{M} \sum_{j=1}^{n} P_{i}^{j}$$
(10)

where S^{j} is the set of activated experts for token 279 j, which is calculated by Equation 6, and P^{j} is the probability of selecting each experts for token *j*, calculated by Equation 2. 282

Final Loss Our model is a generative model that uses next token generation as the training objective. We denote this loss as $Loss_{lm}$. Our final loss is a combination of the language model loss, dynamic loss, and load-balance loss:

$$Loss = Loss_{lm} + \alpha Loss_b + \beta Loss_d \quad (11)$$

where α and β are hyper-parameters to adjust the contribution of the load balance loss and dynamic loss, respectively. In our experiment, we set α as 1e-2 and β is set as 1e-4.

3 **Experiments**

3.1 Settings

3.1.1 **Training data**

We use RedPajama(Computer, 2023) as our training data, which is a fully open-source implementation of the LLaMa dataset. RedPajama data consists of diverse sources including the common crawl (CC), C4, github, Wikipedia, books, arxiv and Stackexchange. In our main experiments, we train all models for 100B tokens.

3.1.2 Model Settings

The model architecture follows LLaMA(Touvron et al., 2023). We use Llama2 tokenizer whose vocabulary size is 32,000. The number of transformer layers is 24 and the hidden dimension is 1024. Each MoE layer has 16 experts. Under this configuration, dense model has approximately 374M parameters. Each MoE model has 3.5B total parameters. Only 374M parameters are activated in MoE-Top1 and 581M parameters are activated in MoE-Top2. More detailed model and training settings are shown in Appendix 9.

3.1.3 Evaluation

We use opencompass¹ to evaluate our model.

3.1.4 Experiment Models

We train several variants of our architecture from scratch using the above model settings.

Dense We use dense models as our baseline. In dense models, each transformer layer is composed of a multi-head attention layer and a standard Feed Forward Network. We implement two Dense models: Dense(374M) and Dense(570M) by setting the hidden dimensions to 1024 and 1280, respectively.

¹https://github.com/open-compass/OpenCompass/

	Dense(374M)	Dense(570M)	MoE-Top1	MoE-Top2	MoE-Dynamic
PIQA (Bisk et al., 2020)	64.3	65.9	67.3	68.1	68.1
Hellaswag (Zellers et al., 2019)	36.1	39.6	42.3	43.9	44.3
ARC-e (Bhakthavatsalam et al., 2021)	37.9	37.6	39.5	40.4	39.9
Commonsense QA (Talmor et al., 2019)	32.2	31.7	30.3	32.1	33.6
BBH (Suzgun et al., 2023)	22.3	22.1	23.0	23.3	25.6
Avg	38.6	39.4	40.5	41.6	42.3

Table 1: Performance on downstream tasks. The best result for each task is emphasized in **bold**.

326MoE-Top1 / Top2The MoE models with Top-K327routing, where K = 1 and 2, respectively. Only328language modeling loss, $Loss_{lm}$, and load-balance329loss $Loss_b$ are used for training. The MoE-Top1330could be seen as a re-implementation of Switch331Transformer (Fedus et al., 2022) and the MoE-Top2332is a re-implementation of Gshard(Lepikhin et al.,3332021). The activated parameters of MoE-Top1 and334MoE-Top2 are nearly the same as Dense(374M)335and Dense(570M), respectively.

MoE-Dynamic MoE-Dynamic model uses our dynamic adaptive routing mechanism, activating a various number of experts depending on the input token representation. The threshold p in our routing mechanism is 0.4. During inference, MoE-Dynamic model activates no more than 2 experts, which means it uses fewer parameters than MoE-Top2.

3.2 Main Results

336

337

339

340

341

342

346

347

351

358

366

Table 1 shows the performance of different models on downstream tasks. Overall, the MoE models outperform the Dense models. Among all the MoE variants, our proposed Dynamic Adaptive MoE demonstrates the best performance, achieving at least a 0.7% higher score on average compared to other models.

We first compare models with an equal number of activated parameters. It is observed that MoE-Top1 outperforms the Dense model with 374M parameters by an average of 1.9% score, and MoE-Top2 surpasses the Dense model with 570M parameters by 2.2% score. This indicates that, with the same number of activated parameters, MoE models substantially outshine their corresponding Dense counterparts.

When comparing models with the same architecture, we generally observe a positive correlation between model performance and the number of activated parameters. For Dense models, the model with 570M parameters outperforms the model with 374M parameters by 0.8% score on average. Sim-



Figure 2: Average scores of MoE-Dynamic with different threshold p on downstream tasks

367

368

369

370

371

373

374

375

376

377

378

379

380

381

383

386

387

390

391

392

393

394

395

396

ilarly, among models using the MoE architecture with a fixed number of activated experts, MoE-Top2, which activates two experts, reaches an average of 41.6% score, outperforming MoE-Top1, which only activates one expert, by 1.1% score. In fact, MoE-Top2 performs better than MoE-Top1 in all subtasks, demonstrating the rule of more parameters leading to better performance.

However, our proposed Dynamic Routing mechanism breaks this rule. As shown in Table 3, the average number of activated experts in the MoE-Dynamic during evaluation phases is less than two, meaning it activates fewer parameters than MoE-Top2. Yet, as shown in Table 1, compared to MoE-Top2, MoE-Dynamic achieves comparable or even better performance on nearly all the tasks and outperforms MoE-Top2 by 0.7% score on average. MoE-Dynamic obtains better performance, indicating that our dynamic routing mechanism can allocate the necessary experts for different inputs more reasonably and make use of parameters more efficiently.

3.3 Effect of Threshold p

Threshold p is a hyper-parameter used to control the dynamic routing mechanism. Training models from scratch with different values of p is resourceintensive. Hence, we explore the impact of this hyper-parameter by performing inference on a pretrained model with varying values of p from 0.1 to 0.7. Table 2 demonstrates the average performance



Figure 3: Average activated experts number across training procedure.

on downstream tasks with different p.

397

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

The table reveals that when p is too low, like 0.1 and 0.2, the model's performance on downstream tasks markedly decreases due to the activation of too few experts. Conversely, once p surpasses a certain threshold, the model's performance stabilizes, and the impact of this parameter on downstream tasks will become minimal.

4 Efficiency of Dynamic Routing

The greatest advantage of MoE models is their ability to efficiently scale to larger models. The Top-K routing mechanism controls the number of parameters used by the entire model by activating a fixed number of experts. In contrast, our proposed dynamic routing mechanism removes the limitation of a fixed number of experts. Naturally, there may be concerns that our method might assign too many experts to each token. To address these concerns, we demonstrate the efficiency of the dynamic routing mechanism from both training and inference perspectives.

4.1 Efficient Training

We sample 1000 pieces of data from different 419 420 sources within Redpajama and calculate the average number of experts activated per token at differ-421 ent stages of training. Figure 3 shows the change in 422 the average number of experts activated throughout 423 the training process of 100B tokens. From the fig-424 ure, we can see that the number of experts activated 425 per token decreases over time. In the early stages 426 of training, dynamic routing assigns more experts 427 to each token, but after 60B tokens, the average 428 number of activated experts is already less than 2. 429 430 Table 2 displays the number of experts activated by MoE-Dynamic at the end of the 100B training. It 431 is evident that across all data sources, the number 432 of experts activated by MoE-Dynamic is less than 433 2. 434

Sources	Ratio	Activated Experts
CC	67%	1.82
C4	15%	1.84
Github	4.5%	1.88
Wiki	4.5%	1.78
Book	4.5%	1.73
Arxiv	2.5%	1.90
StackExchange	2%	1.79
Avg	100%	1.82

Table 2: Average activated experts in different parts of the training corpus.

Recently, the amount of tokens used in training for large language models far exceeds 100B, for instance, Pythia uses 300B tokens, and Llama2 uses 2T tokens. If we continue to train on an even larger scale corpus, the average number of parameters used throughout the training process is guaranteed to be lower than that of Top2-Routing. 435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

4.2 Efficient Inference

To further explore whether our proposed method is efficient in inference, we calculate the average number of experts activated by the model across different downstream tasks. For every question, we use the template from the evaluation to concatenate the question with the gold answer into a complete input and truncate the tokens exceeding 2048 to fit our model's maximum input length. Table 3 shows the average number of experts activated per token across various downstream tasks. The result is averaged across all the layers of transformers and it is evaluated using the checkpoint trained on 100B tokens.

From the table, we can observe that across all five downstream tasks, the number of activated experts is less than two, averaging 1.76 activated experts, which is fewer than the fixed activation of two experts by the Top2 routing method. During the training phase, our method and Top2 routing are comparable in efficiency, but upon completion of training, our inference efficiency substantially outperforms Top2 routing. Given that models are mostly trained once with a greater burden placed on the subsequent deployment nowadays, the advantages of our method over traditional MoE routing mechanisms like Top2 become even more apparent.

5 What is Challenging Input?

The motivation for designing dynamic routing is to enable the model to dynamically adjust the number of allocated experts based on the difficulty of the

6

Sources	Activated Experts
PIQA	1.72
Winogrande	1.76
ARC-e	1.73
Commonsense QA	1.74
BBH	1.87
Avg	1.76

Table 3: Average activated experts in different downstream tasks.

input. In this section, we will explore what kindsof inputs are considered challenging for the modelfrom various perspectives.

5.1 Tasks Requiring Reasoning

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

503

505

507

508

510

511

From Table 3, we could observe that solving the BBH task requires activating an average of 1.87 experts, more than the number needed for other tasks. BBH, which stands for BIG-Bench Hard, is a suite of 23 challenging BIG-Bench tasks. These tasks demand capabilities such as multi-hop reasoning, causal inference, logical deduction, and so on, making them substantially more difficult than normal NLP tasks (Suzgun et al., 2022). Our model's use of more experts on BBH tasks implies that our method indeed can dynamically monitor task difficulty and apply more parameters to tackle more challenging tasks. Interestingly, as shown in Table 1, MoE-Dynamic, compared to MoE-Top2, sees the most improvement on BBH tasks. While the average improvement across all tasks is less than 1.0%, the improvement on BBH is more than 2.0%, which is more than double that of other tasks. This further illustrates that dynamically adjusting the number of activated experts is beneficial for solving downstream tasks, especially more challenging ones.

5.2 Tokens with Ambiguous Semantics

To further analyze what types of tokens are considered more challenging for a model, we examine the average number of experts activated for each token in the vocabulary across different contexts.

We sample 1 million tokens from each part of the training dataset Redpajama, like arxiv and CC, resulting in a new corpus of a total of 7 million tokens. In this corpus, we calculate the average number of experts activated for each token in the vocabulary. To minimize the effect of randomness, we only consider tokens that appear more than 10,000 times in the corpus.

	Examples	C-Words Ratio
Most Experts	tr, eq, mu, frac	10
Least Expers	to, that, and, show	51

Table 4: The first column shows examples of tokens requiring the most experts and least experts. The last column shows the complete words ratio in these two groups of tokens.

Table 4 shows the number of complete words among the top 100 and bottom 100 tokens by the average number of experts activated, along with some examples. 512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

Upon manually reviewing the 100 tokens that activate the most experts and the 100 tokens that activate the least, we observe an interesting phenomenon: Tokens with relatively definite semantics are considered easier by the model, activating fewer experts. In contrast, tokens with uncertain semantics are deemed more challenging and require more experts for processing.

Specifically, since our model's tokenizer is trained with Byte Pair Encoding (BPE), many tokens are not complete words but subwords. These subwords have vaguer semantics compared to full words because they can combine with many other subwords to form words with different meanings. For example, the subword 'tr' can lead to the formation of hundreds of words with varied meanings, such as tree, triple, train, trick, trouble, and so on. Due to the multitude of possible semantics, different meanings may require different experts for processing, making such subwords require a comprehensive understanding by more experts.

6 Bottom Layers Need More Experts

An intriguing observation from our study is that our model achieves superior performance while activating fewer parameters. As shown in Table 3, on all the tasks, our MoE-Dynamic activates an average of fewer than two experts. But it outperforms the MoE-Top2 in downstream tasks as shown in Table 1. This result is quite surprising, as performance on downstream tasks is typically correlated with the quantity of activated parameters. We attribute this unexpected phenomenon to our method's more proper allocation of the experts to be activated across different layers, employing more experts at lower levels and fewer at the top. This layer-wise dynamic allocation, as opposed to the fixed number of experts per layer, somewhat mitigates the common issue of overthinking in deep



Figure 4: Activated experts in different layers

neural networks, thereby enhancing performance.

The overthinking refers to the situations where simpler representations of an input sample at an earlier layer, relative to the complex representations at the final layer, are adequate to make a correct prediction (Kaya et al., 2019). Previous works (Liu et al., 2020; Schwartz et al., 2020; Xin et al., 2021) have demonstrated that shallower representations can achieve comparable, if not better, performance across various tasks than deeper representations. This could be due to deeper representations overfitting specific distributions, lacking generalizability, and being more vulnerable to attacks (Hu et al., 2019; Zhou et al., 2020). It suggests that in some cases, acquiring a better shallow representation is more valuable than obtaining a more complex deep representation, which correlates to previous findings that removing top layers has a limited impact on the downstream tasks (Sajjad et al., 2023).

Compared with Top2 routing, our dynamic adaptive routing activates more experts at the bottom layers to obtain better shallow representations and use the simpler network in the top layers to alleviate the overthinking issue. Figure 4 displays the number of experts activated per token at different layers². From the figure, we observe a gradual decrease in the average number of experts activated per token with increasing layer depth. The lowest layer activates the most experts, up to 4 experts per token, enabling better shallow representations through a wider network, which is beneficial for various downstream tasks. At the topmost layer, the number of activated experts per token is reduced to even one. This phenomenon can avoid model being too complex and preserve generality in the final representation.

7 Related Work

The Mixture of Experts (MoE) model is initially introduced by (Jacobs et al., 1991). Recent studies have demonstrated sparsely gated MoE models have substantial improvements in model capacity and efficiency, enabling superiors performance than dense models(Shazeer et al., 2017). Particularly MoE has shown great potential with the integration of transformer architectures (Zoph et al., 2022). 590

591

592

593

594

595

597

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

In previous MoE architectures, a static number of experts are activated regardless of the varying complexity presented by input tokens. Most of MoE models activate Top-1 or Top-2 experts (Lepikhin et al., 2021; Fedus et al., 2022), which could potentially limit the efficacy of MoE models.

There are works allocating various number of experts for input tokens. Expert-Choice MoE model selects Top-K tokens for each expert(Zhou et al., 2022). However, in Expert-choice MoE model, the floating-point operations per second(FLOPS) in each MoE layer are the same. Previous work indicates that different MoE layers may need different FLOPS to achieve optimal performance(Jawahar et al., 2023).

Different from these prior works, our dynamic routing mechanism can allocate more experts for complex tokens and fewer for simpler ones. Additionally, it strategically selects more experts in the lower layers and fewer in the upper layers, thereby minimizing computational redundancy. Experimental results demonstrate that this dynamic routing approach contributes to improvements in both the efficiency and performance of MoE models.

8 Conclusion

Our paper introduces a dynamic expert selection framework for Mixture of Experts (MoE) models, surpassing traditional fixed Top-K routing by adjusting expert activation based on input complexity. Our approach not only improves computational efficiency but also model performance, evidenced by obvious gains over conventional Top-K routing in our evaluations. Our findings reveal the framework's effectiveness at dynamically dispatching different numbers of experts, particularly for complex reasoning tasks, and suggest the potential for developing more challenging heterogeneous MoE models. In support of further research, we will open-source our models, contributing to advancements in the MoE domain.

589

554

²The results are evaluated using a checkpoint trained on 100B tokens.

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

693

694

Limitation

639

655

656

657

665

674

675

676

677

678

Due to resource constraints, the size of the model we trained is limited, with only about 600M acti-641 vation parameters, and the entire MoE (Mixture of Experts) model being just over 3B in size. How-643 ever, (Dai et al., 2024) has validated that within the MoE framework, conclusions drawn from smaller models can be generalized to larger models with more parameters. Hence, we believe our proposed 647 dynamic routing method could also be effective in larger-scale models. Additionally, we have only trained on 100B tokens, which may not be sufficient for model training. Yet, given the same scale of training data, our method demonstrated superior performance, which also underscores the efficiency of our training process. 654

References

- Sumithra Bhakthavatsalam, Daniel Khashabi, Tushar Khot, Bhavana Dalvi Mishra, Kyle Richardson, Ashish Sabharwal, Carissa Schoenick, Oyvind Tafjord, and Peter Clark. 2021. Think you have solved direct-answer question answering? try arc-da, the direct-answer AI2 reasoning challenge. *CoRR*, abs/2102.03315.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. PIQA: reasoning about physical commonsense in natural language. In The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, pages 7432– 7439. AAAI Press.
- Together Computer. 2023. Redpajama: An open source recipe to reproduce llama training dataset.
- 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.
- Nan Du, Yanping Huang, Andrew M. Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten P. Bosma, Zongwei Zhou, Tao Wang, Yu Emma Wang, Kellie Webster, Marie Pellat, Kevin Robinson, Kathleen S. Meier-Hellstern, Toju Duke, Lucas Dixon, Kun Zhang, Quoc V. Le, Yonghui Wu, Zhifeng Chen, and Claire Cui. 2022. Glam: Efficient scaling of language models with mixture-of-experts. In International Conference on Machine Learning, ICML 2022, 17-23

July 2022, Baltimore, Maryland, USA, volume 162 of Proceedings of Machine Learning Research, pages 5547–5569. PMLR.

- William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *J. Mach. Learn. Res.*, 23:120:1–120:39.
- Ting-Kuei Hu, Tianlong Chen, Haotao Wang, and Zhangyang Wang. 2019. Triple wins: Boosting accuracy, robustness and efficiency together by enabling input-adaptive inference. In *International Conference on Learning Representations*.
- Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. 1991. Adaptive mixtures of local experts. *Neural Comput.*, 3(1):79–87.
- Ganesh Jawahar, Subhabrata Mukherjee, Xiaodong Liu, Young Jin Kim, Muhammad Abdul-Mageed, Laks V. S. Lakshmanan, Ahmed Hassan Awadallah, Sébastien Bubeck, and Jianfeng Gao. 2023. Automoe: Heterogeneous mixture-of-experts with adaptive computation for efficient neural machine translation. In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 9116–9132. Association for Computational Linguistics.
- Yigitcan Kaya, Sanghyun Hong, and Tudor Dumitras. 2019. Shallow-deep networks: Understanding and mitigating network overthinking. In *International conference on machine learning*, pages 3301–3310. PMLR.
- Young Jin Kim, Ammar Ahmad Awan, Alexandre Muzio, Andrés Felipe Cruz-Salinas, Liyang Lu, Amr Hendy, Samyam Rajbhandari, Yuxiong He, and Hany Hassan Awadalla. 2021. Scalable and efficient moe training for multitask multilingual models. *CoRR*, abs/2109.10465.
- Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. 2021. Gshard: Scaling giant models with conditional computation and automatic sharding. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net.
- Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. 2021. BASE layers: Simplifying training of large, sparse models. In Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event, volume 139 of Proceedings of Machine Learning Research, pages 6265–6274. PMLR.
- Weijie Liu, Peng Zhou, Zhiruo Wang, Zhe Zhao, Haotang Deng, and Qi Ju. 2020. Fastbert: a selfdistilling bert with adaptive inference time. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6035– 6044.

- 750 751
- 75
- 754 755
- 757 758 759
- 7
- 762 763
- 764 765
- 7
- 7
- 769 770 771
- 772 773 774
- 775
- 777
- 778 779
- 7777
- 7
- 787 788 789

- 791 792 793
- 794
- 795 796
- 797 798
- 801 802
- 8
- 8

805

- Hassan Sajjad, Fahim Dalvi, Nadir Durrani, and Preslav Nakov. 2023. On the effect of dropping layers of pre-trained transformer models. *Computer Speech & Language*, 77:101429.
- Roy Schwartz, Gabriel Stanovsky, Swabha Swayamdipta, Jesse Dodge, and Noah A Smith. 2020. The right tool for the job: Matching model and instance complexities. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6640–6651.
- Noam Shazeer. 2020. GLU variants improve transformer. *CoRR*, abs/2002.05202.
 - Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc V. Le, Geoffrey E. Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V. Le, Ed H. Chi, Denny Zhou, and Jason Wei. 2023. Challenging big-bench tasks and whether chain-of-thought can solve them. In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 13003–13051. Association for Computational Linguistics.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. *arXiv preprint arXiv:2210.09261*.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 4149–4158. Association for Computational Linguistics.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971.
- Ji Xin, Raphael Tang, Yaoliang Yu, and Jimmy Lin. 2021. Berxit: Early exiting for bert with better finetuning and extension to regression. In *Proceedings* of the 16th conference of the European chapter of the association for computational linguistics: Main Volume, pages 91–104.

Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In *Proceedings* of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers, pages 4791–4800. Association for Computational Linguistics. 807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

- Wangchunshu Zhou, Canwen Xu, Tao Ge, Julian McAuley, Ke Xu, and Furu Wei. 2020. Bert loses patience: Fast and robust inference with early exit. *Advances in Neural Information Processing Systems*, 33:18330–18341.
- Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M. Dai, Zhifeng Chen, Quoc V. Le, and James Laudon. 2022. Mixture-ofexperts with expert choice routing. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.
- Barret Zoph, Irwan Bello, Sameer Kumar, Nan Du, Yanping Huang, Jeff Dean, Noam Shazeer, and William Fedus. 2022. St-moe: Designing stable and transferable sparse expert models. *arXiv preprint arXiv:2202.08906*.
- Simiao Zuo, Xiaodong Liu, Jian Jiao, Young Jin Kim, Hany Hassan, Ruofei Zhang, Jianfeng Gao, and Tuo Zhao. 2022. Taming sparsely activated transformer with stochastic experts. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net.

9 Detailed Training Setting

9.1 Model Setting

The model architecture follows LLaMA(Touvron 841 et al., 2023). We use Llama2 tokenizer whose vo-842 cabulary size is 32000. Unless specifically stated 843 otherwise, we set the number of transformer layers 844 to 24, the hidden dimension to 1024. We employ 845 the multi-head attention mechanism with a total of 846 16 attention heads, where each head has a dimen-847 sion of 64. We use SwiGLU(Shazeer, 2020) in FFN 848 layers. For initialization, all learnable parameters 849 are randomly initialized with a standard deviation 850 of 0.006. Each MoE layer has 16 experts, which 851 have the same initialized parameters as a standard 852 FFN. Under this configuration, each dense model 853 has has approximately 374M parameters. Each 854 MoE model has 3.5B total parameters. Only 374 855 parameters are activated in MoE-Top1 and 581M 856 parameters are activated in MoE-Top2. 857

9.2 Training Setting

858

We adopt AdamW optimizer with first-moment decay $\beta_1 = 0.9$ and second-moment decay $\beta_2 =$ 0.95. The weight decay is 0.1. The learning rate warms up from 0 to 3e-4 in the first 2000 steps and decays in the remaining steps using the cosine decay schedule to 3e-5. We set the context length to 2048 and adopt the batch size of 2048.