

# On the Effectiveness of Discrete Representations in Sparse Mixture of Experts

Anonymous ACL submission

## Abstract

Sparse Mixture of Experts (SMoE) is an effective solution for scaling up model capacity without increasing the computational costs. A crucial component of SMoE is the router, responsible for directing the input to relevant experts; however, it also presents a major weakness, leading to routing inconsistencies and representation collapse issues. Instead of fixing the router like previous works, we propose an alternative that assigns experts to input via *indirection*, which employs the discrete representation of input that points to the expert. The discrete representations are learnt via vector quantization, resulting in a new architecture dubbed Vector-Quantized Mixture of Experts (VQMoE). We provide theoretical support and empirical evidence demonstrating the VQMoE’s ability to overcome the challenges present in traditional routers. Through extensive evaluations on both large language models and vision tasks for pre-training and fine-tuning, we show that VQMoE achieves a 28% improvement in robustness compared to other SMoE routing methods while maintaining strong performance in fine-tuning tasks.

## 1 Introduction

Scaling Transformers with data and compute has demonstrated unprecedented successes across various domains such as natural language processing (NLP) tasks (Du et al., 2022; Fedus et al., 2022; Zhou et al., 2024), and visual representation learning (Riquelme et al., 2021a; Shen et al., 2023b).

However, training and inference of a single large Transformer-based model might require hundreds of thousands of compute hours, costing millions of dollars (Kaddour et al., 2023). This issue has motivated contemporary studies to investigate Sparse Mixture of Experts (SMoE) (Shazeer et al., 2017; Zoph et al., 2022; Xue et al., 2024; Jiang et al., 2024). SMoE models that are inspired by (Jacobs

et al., 1991a) usually include a set of experts sharing the same architecture and a router that activates only one or a few experts for each input. Compared to dense models of the same size, SMoE counterparts significantly reduce inference time thanks to not using all experts simultaneously (Artetxe et al., 2022; Krajewski et al., 2024).

However, training SMoEs remains a challenge due to representation collapse, that is, either a small number of experts receive most of the routed tokens or all experts converge to learn similar representations. To tackle the issue, several works (Chi et al., 2022; Chen et al., 2023a; Do et al., 2023) have focused on router policy improvement. However, these do not touch a fundamental question, ‘Do we really need a router in the first place?’ Our research suggests that adopting a discrete representation could help solve the challenges currently faced by the router method. Discrete representation learning in the context of SMoE is motivated by its ability to capture structured and interpretable patterns within data, aligning with the way that humans categorize and process information through distinct symbols, like tokens. This approach enables better generalization and facilitates knowledge transfer across different contexts. Additionally, discrete representations provide a robust and efficient mechanism for selecting and routing inputs to the appropriate experts by clustering them more effectively. By bridging the gap between discrete and continuous representations, this method leads to more stable and interpretable expert assignments, helping to mitigate issues such as representation collapse and overfitting, which are common challenges in SMoE training.

Employing vector quantization (VQ) techniques to learn discrete representation, this paper proposes a novel mixture of expert framework, named VQMoE, which overcomes the representation collapse and inconsistency in training sparse mixture of experts. More specifically, we prove that the exist-

ing router methods are inconsistent and VQMoE suggests an optimal expert selection for training SMoE. Additionally, our method guarantees superior SMoE training strategies compared to the existing methods by solving the representation collapse by design.

We evaluate the proposed method by conducting pre-training of Large Language Models (LLMs) on several advanced SMoE architectures, such as SMoE (Jiang et al., 2024), StableMoE (Dai et al., 2022), or XMoE (Chi et al., 2022), followed by fine-tuning on downstream tasks on both Language and Vision domains.

In summary, the primary contributions of this paper are threefold: (1) we theoretically demonstrate that learning a discrete representation is an optimal approach for expert selection and that VQMoE inherently addresses the issue of representation collapse; (2) we propose the use of the Vector Quantization method to learn cluster structures and resolve related challenges; and (3) we conduct extensive experiments on large language models and vision pre-training and fine-tuning tasks, providing an in-depth analysis of VQMoE’s behavior to showcase its effectiveness.

## 2 Related Work

**Sparse Mixture of Experts (SMoE).** Sparse Mixture of Experts (SMoE) builds on the Mixture of Experts (MoE) framework introduced by Jacobs et al. (1991b); Jordan and Jacobs (1994), with the core idea that only a subset of parameters is utilized to process each example. This approach was first popularized by Shazeer et al. (2017). SMoE’s popularity surged when it was combined with large language models based on Transformers (Zhou et al., 2022; Li et al., 2022; Shen et al., 2023a), and its success in natural language processing led to its application across various fields, such as computer vision (Riquelme et al., 2021b; Hwang et al., 2023; Lin et al., 2024), speech recognition (Wang et al., 2023; Kwon and Chung, 2023), and multi-task learning (Ye and Xu, 2023; Chen et al., 2023b).

However, SMoE faces a major problem in training known as representation collapse, i.e., the experts converge to similar outputs. To address this, various methods have been introduced. XMoE (Chi et al., 2022) calculates routing scores between tokens and experts on a low-dimensional hypersphere. SMoE-dropout (Chen et al., 2023a) uses a fixed, randomly initialized router network to activate experts and gradually increase the number

of experts involved to mitigate collapse. Similarly, HyperRouter (Do et al., 2023) utilizes HyperNetworks (Ha et al., 2016) to generate router weights, providing another pathway for training SMoE effectively. StableMoE (Dai et al., 2022) introduces a balanced routing approach where a lightweight router, decoupled from the backbone model, is distilled to manage token-to-expert assignments. The StableMoE strategy ensures stable routing by freezing the assignments during training, while SimSMoE (Do et al., 2024) forces experts to learn dissimilar representations. Despite these extensive efforts, the representation collapse issue persists, as highlighted by Pham et al. (2024). While most solutions focus on improving routing algorithms, our approach takes a different path by learning a discrete representation of input that points to relevant experts.

**Discrete Representation.** Discrete representations align well with human thought processes; for example, language can be understood as a series of distinct symbols. Nevertheless, the use of discrete variables in deep learning has proven challenging, as evidenced by the widespread preference for continuous latent variables in most current research. VQVAE (van den Oord et al., 2017) implements discrete representation in Variational AutoEncoder (VAE) (Kingma and Welling, 2022) using vector quantisation (VQ). IMSAT (Hu et al., 2017) attains a discrete representation by maximizing the information-theoretic dependency between data and their predicted discrete representations. Recent works follow up the vector quantisation ideas and make some enhancements for VAE, for example: (Yu et al., 2022); (Mentzer et al., 2023); and (Yang et al., 2023). Mao et al. (2022) utilize a discrete representation to strengthen Vision Transformer (ViT) (Dosovitskiy et al., 2021). To the best of our knowledge, our paper is the first to learn a discrete representation of Sparse Mixture of Experts.

## 3 Method

We propose a novel model, Vector-Quantized Mixture of Experts (VQMoE), which learns discrete representations for expert selection. As illustrated in Fig. 1a, our approach selects experts directly based on the input representation, eliminating the need for a trained router. To prevent information loss, we integrate discrete and continuous representations within the model.

### 3.1 Preliminaries

**Sparse Mixture of Experts.** Sparse Mixture of Experts (SMoE) is often a transformer architecture that replaces the MLP layers in standard transformers with Mixture of Experts (MoE) layers (Shazeer et al., 2017). Given  $\mathbf{x} \in \mathbb{R}^{n \times d}$  as the output of the multi-head attention (MHA), the output of SMoE with  $N$  experts is a weighted sum of each expert’s computation  $E_i(\mathbf{x})$  by the router function  $\mathcal{S}(\mathbf{x})$ :

$$\begin{aligned} f_{\text{SMoE}}(\mathbf{x}) &= \sum_{i=1}^N \mathcal{S}(\mathbf{x})_i \cdot E_i(\mathbf{x}) \\ &= \sum_{i=1}^N \mathcal{S}(\mathbf{x})_i \cdot \mathbf{W}_{\text{FFN}_i}^2 \phi(\mathbf{W}_{\text{FFN}_i}^1 \mathbf{x}) \end{aligned} \quad (1)$$

Where  $\mathcal{S}(\mathbf{x})$  is computed by *TopK* function as equation (2) that determines the contribution of each expert to the SMoE output.

$$\begin{aligned} \mathcal{S}(\mathbf{x}) &= \text{TopK}(\text{softmax}(\mathcal{G}(\mathbf{x})), k), \\ \text{TopK}(\mathbf{v}, k) &= \begin{cases} v_i & \text{if } v_i \in \text{top } k \text{ largest of } \mathbf{v}, \\ -\infty & \text{otherwise.} \end{cases} \end{aligned} \quad (2)$$

**Discrete Representation Learning.** van den Oord et al. (2017) propose VQVAE, which uses Vector Quantisation (VQ) to learn a discrete representation. Given an input  $x \in \mathbb{R}^{n \times d}$ , VQVAE discretized the input into a codebook  $V \in \mathbb{R}^{K \times d}$  where  $K$  is the codebook size and  $d$  is the dimension of the embedding. Let denote  $z_v(x) \in \mathbb{R}^{n \times d}$  denotes the output of the VQVAE and  $\mathbf{1}(\cdot)$  is the indicator function. The discrete representation  $z_q(x_i) = v_k$ , where  $k = \text{argmin}_j \|z_v(x_i) - v_j\|_2$  is achieved by vector quantizer  $q_\theta$  that maps an integer  $z$  for each input  $x$  as:

$$q_\theta(z = k \mid x) = \mathbf{1} \left( k = \arg \min_{j=1:K} \|z_v(x) - V_j\|_2 \right) \quad (3)$$

### 3.2 Vector-Quantized Mixture of Experts (VQMoE)

**Pre-training VQMoE.** Existing Sparse Mixture of Experts (SMoE) models learn continuous representations and select experts based on routing scores derived from token-expert embeddings. In this paper, we propose a novel architecture that learns simultaneously continuous and discrete representations at a training phase as Figure 1a. The

continuous representation enables the model to capture complex structures in the data, while the discrete representation learns latent representation from data and then transfers the knowledge to downstream tasks. Given  $\mathbf{x} \in \mathbb{R}^{n \times d}$  as the output of the MHA and  $f^v$  is a vector quantization operator, the output of the VQMoE layer at the Pre-training phase as follows:

$$f^{\text{VQMoE}}(\mathbf{x}) = g(\mathbf{x})_c f^{\text{SMoE}}(\mathbf{x}) + g(\mathbf{x})_d \sum_{l=1}^K f_l^{\text{FFN}}(\tilde{\mathbf{x}}_l) \quad (4)$$

Where  $\tilde{\mathbf{x}}_l = v_k$  if  $x_l \in V_l$  codebook, otherwise  $\tilde{\mathbf{x}}_l = \vec{0}$ ;  $f_l^{\text{FFN}}(\tilde{\mathbf{x}}_l)$  corresponds to the expert associated with the  $V_l$  codebook;  $g(\mathbf{x})_c(\mathbf{x}) = \text{col}_0(G(\mathbf{x}))$ ,  $g(\mathbf{x})_d(\mathbf{x}) = \text{col}_1(G(\mathbf{x}))$  is gating function for continuous and discrete representation with  $G(\mathbf{x}) = \text{softmax}(\mathbf{W}_g^T \times \mathbf{x})$ .  $\mathbf{W}_g^T \in \mathbb{R}^{2 \times d}$  is a learnable weight and  $K$  is number of codes.

**Fine-tuning VQMoE.** According to (Geva et al., 2021), the Feed-forward layers (FFN) constitute two-thirds of a transformer model’s parameters. Thus, VQMoE enhances the robustness and efficiency of the Mixture of Experts by leveraging the discrete representations learned during the Pre-training phase. For further details, the output of VQMoE during the fine-tuning stages only requires the discrete representation part as Figure 1b, leading to the following output from the VQMoE layer in the fine-tuning phase:

$$f^{\text{VQMoE}}(\mathbf{x}) = \sum_{l=1}^K f_l^{\text{FFN}}(\tilde{\mathbf{x}}_l) \quad (5)$$

### 3.3 Training Procedure

**Pretraining.** The training objective is jointly minimizing the loss of the target task and losses of the Vector Quantization module ( $\mathcal{L}^{\text{l2}}$  and  $\mathcal{L}^{\text{commitment}}$ ) as in (van den Oord et al., 2017). Equation 6 specifies the overall loss function for training VQMoE with three components: (1) task loss; (2)  $l_2$  loss; (3) a commitment loss. While  $\mathcal{L}^{\text{l2}}$  helps to move the embedding  $v_i$  towards the outputs  $z_v(x)$ , the commitment loss makes sure the output of the Vector Quantization module commits to the embedding and its output does not grow. The Vector Quantization algorithm does not vary with  $\beta$ , we follow  $\beta = 0.25$  as van den Oord et al. (2017). We introduce a new parameter,  $\alpha$ , to regulate the contribution of the Vector Quantization loss to the overall loss. A higher value of  $\alpha$  favors a stronger adherence to the discrete representation, and vice versa.

$$L = \mathcal{L}_{\text{task}} + \alpha(\|\text{sg}[z_v(x)] - v\|_2^2 + \beta\|z_v(x) - \text{sg}[v]\|_2^2) \quad (6)$$

where  $\text{sg}(\cdot)$  is the stop gradient operator defined as follows:

$$\text{sg}(x) = \begin{cases} x & \text{forward pass} \\ 0 & \text{backward pass} \end{cases} \quad (7)$$

**Fine-tuning.** For downstream tasks, we fine-tune the pretraining model by utilizing the codebook learned from the Equation 6 by freezing all parameters at the Vector Quantization module. Thus, the training objective simply becomes:  $L = \mathcal{L}_{\text{task}}$ .

#### 4 VQMoE solves Representation Collapse by Design

The representation collapse problems in SMoE, which leads all experts to learn the same thing, first declared by (Chi et al., 2022). Same as (Chi et al., 2022); (Do et al., 2023), we illustrate the presentation collapse issue by the Jacobian matrix approach. Indeed, Jacobian matrix of SMoE with respect to  $x \in \mathbb{R}^{n \times d}$  is followed as:

$$\begin{aligned} J_{\text{SMoE}} &= \mathcal{S}(x)_k \mathbf{J}^{\text{FFN}} + \sum_{j=1}^N \mathcal{S}(x)_k (\delta_{kj} - S_j) \mathbf{E}(x)_i e_j^\top \\ &= \mathcal{S}(x)_k \mathbf{J}^{\text{FFN}} + \sum_{j=1}^N c_j e_j^\top. \end{aligned} \quad (8)$$

where  $c_j = \mathcal{S}(x)_k (\delta_{kj} - S_j) \mathbf{E}(x)_i$ . Equation 8 consists two terms: (1)  $\mathcal{S}(x)_k \mathbf{J}^{\text{FFN}}$  represents a contribution from input token and experts to the final output; (2)  $\sum_{j=1}^N c_j e_j^\top$  indicates to learn better gating function to minimize the task loss. Moreover, Equation 8 is suggested to be updated toward a linear combination of the expert embeddings. Since  $N \ll d$  in practice, the above equation shows representation collapse from  $\mathbb{R}^d$  to  $\mathbb{R}^N$ .

Compared to SMoE, does VQMoE reduce the representation collapse issue? To answer the essential question, we calculate the Jacobian matrix of VQMoE with respect to  $x \in \mathbb{R}^{n \times d}$  is given by:

$$\begin{aligned} J_{\text{VQMoE}} &= g(x)_c J_{\text{SMoE}} + J_{g(x)_c} f_{\text{SMoE}}(x) + \\ &\quad g(x)_d J_{\text{VQ}} + J_{g(x)_d} f_{\text{VQMoE}}(x) \end{aligned} \quad (9)$$

Equation 9 is written shortly as below:

$$\begin{aligned} J_{\text{VQMoE}} &= J_1 + \sum_{j=1}^N c_j e_j^\top + \sum_{l=1}^K d_l e_l^\top + \sum_{m \in \{c,d\}} g_m e_m^\top \\ &= J_1 + \sum_{j=1}^{N+K+2} o_j e_j^\top. \end{aligned} \quad (10)$$

where  $J_1 = \mathcal{S}(x)_k \mathbf{J}^{\text{FFN}}$ ;  $c_j = \mathcal{S}(x)_k (\delta_{kj} - S_j) \mathbf{E}(x)_i$ ;  $d_l = g(x)_d$  (due to the vector quantization operator using pass gradient trick (van den Oord et al., 2017));  $g_m = \mathcal{S}(x)_m (\delta_{mk} - S_k) f_m$  where  $f_m \in [f_{\text{SMoE}}(x), f_{\text{VQMoE}}]$ .

Same as the Jacobian matrix of SMoE, the Jacobian matrix of VQMoE consists two terms: (1)  $J_1$  depends on input token and experts to the final output; (2)  $\sum_{j=1}^{N+K+2} o_j e_j^\top$  indicates to learn better gating function to minimize the task loss. We can see that  $N + K + 2 \gg N$ , implying that VQMoE is better than SMoE in solving the representation collapse issue. In theory, we can choose the number of codes to be approximately  $d - N - 2$  with a hashing index to experts to address the issue. However, this involves a trade-off with the computational resources required to learn the codebook.

### 5 Experiment

We conduct experiments to explore the following hypotheses: (i) VQMoE provides an effective SMoE training algorithm for LLMs; (ii) VQMoE delivers a robust and efficient solution during the fine-tuning phase; and (iii) VQMoE outperforms other routing methods in vision domain.

#### 5.1 Experimental Settings

To answer the three above hypotheses, we conduct experiments on Vision, Language, and Time-series tasks. For **Pre-training language models**, we examine two common tasks in the training and evaluation of large language models: character-level language modeling using the enwik8 and text8 datasets (Mahoney, 2011), and word-level language modeling with the WikiText-103 (Merity et al., 2016) and One Billion Word datasets (Chelba et al., 2014). For **Parameter-efficient fine-tuning**, we consider pre-trained base models on enwik8 and efficient Fine-tuning it on a downstream task. We choose the SST-2 (Socher et al., 2013), SST-5 (Socher et al., 2013), IMDB (Maas et al., 2011), and BANKING77 (Casanueva et al., 2020) datasets. For **vision tasks**, we employ the Vision Transformer model (Dosovitskiy et al., 2021) with the



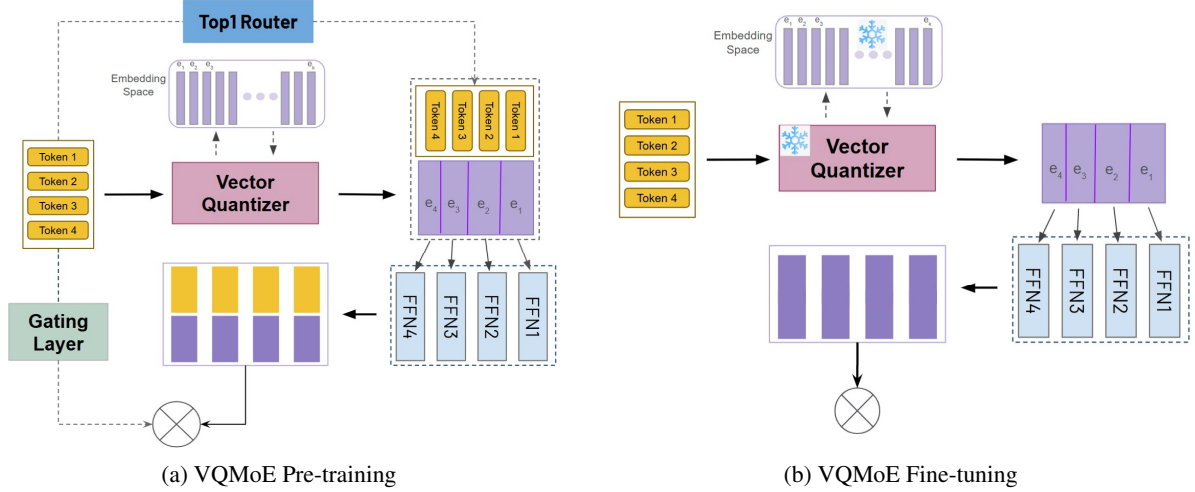


Figure 1: Illustration of the proposed VQMoE architecture for Pre-training and fine-tuning. (a) At the Pre-training stage, VQMoE architecture learns simultaneously continuous and discrete representation at the Pre-training phase. The continuous representation is learned by the conventional SMoE, while the Vector Quantization block facilitates the learning of a discrete representation. The final output is then combined by a gate layer. (b) VQMoE learns a discrete representation that is capable of operating efficiently and robustly on downstream tasks. VQMoE computes the discrete representation only during the fine-tuning stage to achieve robustness and efficiency.

state-of-the-art routing method and our method to train and evaluate the image classification task. Our experiments encompass five widely recognized image classification datasets: Cifar10, Cifar100 (Krizhevsky, 2009), STL-10 (Coates et al., 2011), SVHN (Netzer et al., 2011), ImageNet-1K (Deng et al., 2009).

## 5.2 Pre-training Language Models

**Training tasks** We explore two common tasks in the training and evaluation of LLMs. First, character-level language modeling on the enwik8 or text8 datasets (Mahoney, 2011), which are common datasets to evaluate the model’s pre-training capabilities. We also consider the word-level language modeling task on WikiText-103 (Merity et al., 2016) and One Billion Word dataset (Chelba et al., 2014), a much larger and more challenging dataset, to test the models scaling capabilities. For all datasets, we follow the default splits of training-validation-testing. Second, we consider Fine-tuning the models on downstream applications to investigate the models’ capabilities of adapting to different domains. To this end, we consider pre-trained medium models on enwik8 and Fine-tuning them on a downstream task. We choose the SST-2 (Socher et al., 2013), SST-5 (Socher et al., 2013), IMDB (Maas et al., 2011), and BANKING77 (Casanueva et al., 2020) datasets, which are common NLP tasks to evaluate pre-trained models. Following Chen et al. (2023a), we freeze the router and only optimize the experts’ parameter in this

experiment.

**Models.** For the language tasks, we follow the same settings as in SMoE-Dropout (Chen et al., 2023a). We consider two decoder-only architectures: (i) the standard Transformer (Vaswani et al., 2017); and (ii) and Transformer-XL (Dai et al., 2019a) with the same number of parameters as Transformer. We evaluate our method versus the state of art Sparse Mixture of Expert Layers such as StableMoE (Dai et al., 2022) and XMoe (Chi et al., 2022). We consider two model configurations: (i) base: with four SMoE blocks and **20M** parameters; (ii) large: with twelve SMoE layers and **210M** parameters. We emphasize that we are not trying to achieve state-of-the-art results due to the limited resource constraints. Instead, we evaluate the small and large models on various datasets to demonstrate the scalability and efficacy of our algorithm. Lastly, we conduct extensive investigations using the tiny model to understand the algorithm behaviours and their robustness to different design choices. Lastly, unless otherwise stated, we implement them with  $K = 2$  in the experiments.

**Baselines.** We compare our VQMoE with state-of-the-art SMoE training strategies for LLMs. **SMoE** (Jiang et al., 2024) employs a simple router trained end-to-end with the experts. **StableMoE** (Dai et al., 2022) proposes a two-phase training process where the first phase trains only the router, and then the router is fixed to train the experts in the second phase. **XMoe** (Chi et al.,

Configuration		Enwik8 (BPC)		Text8 (BPC)		WikiText-103 (PPL)		lm1b (PPL)	
Architecture	Algorithm	Base	Large	Base	Large	Base	Large	Base	Large
Transformer	VQMoE	<b>1.48</b>	<b>1.41</b>	<b>1.47</b>	<b>1.40</b>	<b>38.74</b>	<b>31.98</b>	<b>59.48</b>	<b>49.30</b>
	SMoE	1.49	1.41	1.49	1.40	39.50	32.30	60.88	51.30
	SMoE-Dropout	1.82	2.22	1.70	1.89	72.62	107.18	97.45	159.09
	XMoE	1.51	1.42	1.49	1.42	39.56	32.65	61.17	51.84
	StableMoE	1.49	1.42	1.49	1.41	39.45	32.34	60.72	50.74
Transformer-XL	VQMoE	<b>1.19</b>	<b>1.08</b>	<b>1.28</b>	<b>1.17</b>	<b>29.48</b>	<b>23.85</b>	<b>56.85</b>	<b>48.70</b>
	SMoE	1.20	1.09	1.29	1.18	30.16	24.02	58.00	48.71
	SMoE-Dropout	1.56	2.24	1.56	1.86	58.37	40.02	93.17	68.65
	XMoE	1.21	1.09	1.28	1.17	30.34	24.22	58.33	50.64
	StableMoE	1.20	1.10	1.28	1.19	29.97	24.19	58.25	49.17
# Params		20M	210M	20M	210M	20M	210M	20M	210M

Table 1: BPC on the enwik-8 and text8 test sets; and perplexity on the Wikitext-103 and One Billion Word test sets. Lower is better, best results are in bold.

2022) implements a deep router that comprises a down-projection and normalization layer and a gating network with learnable temperatures. Lastly, motivated by SMoE-Dropout (Chen et al., 2023a), we implement the **SMoE-Dropout** strategy that employs a randomly initialized router and freeze it throughout the training process.

**Training procedure.** For the language modeling experiments, we optimize the base models and the large models for 100,000 steps. We use an Adam (Kingma and Ba, 2017) optimizer with a Cosine Annealing learning rate schedule (Loshchilov and Hutter, 2017). The lowest validation loss checkpoint is used to report the final performance on the test set.

**Q1: Does VQMoE perform better on Pre-training tasks compared to routing methods? A1: Yes.**

Table 1 presents the evaluation metrics comparing VQMoE with state-of-the-art approaches. We also show the performance progression of the base model on the validation set. Notably, across all methods, the Transformer-XL architecture consistently outperforms the standard Transformer on all datasets. While advanced strategies like XMoE and StableMoE tend to surpass vanilla SMoE when model complexity is increased (from small to medium) or more data is introduced (moving from enwik8 to WikiText-103 or One Billion Word), these improvements are often inconsistent or marginal. In contrast, VQMoE consistently outperforms all competitors across benchmarks (keeping in mind that the BPC metric is log-scaled), ar-

chitectures, and also converges more quickly. This highlights VQMoE’s effectiveness in learning an efficient routing policy for the language modeling pre-training task.

**Q2: Does VQMoE keep outperforming the router method when scaling up? A2: Yes.**

Table 1 also demonstrates that VQMoE maintains consistently strong performance when scaled up to 12-layer Transformer and Transformer-XL architectures. Across all four datasets, the performance gap between VQMoE and other routing methods widens as the dataset size increases, from enwik8 to the One Billion Word dataset. This suggests that our approach has the potential to scale effectively with larger language models and bigger datasets. An interesting observation is that SMoE-Dropout (Chen et al., 2023a) performs the worst among all methods, indicating that a random routing policy is insufficient and requires updating for effective training. This finding highlights that the success of SMoE-Dropout is largely due to its self-slimmable strategy, which linearly increases the number of activated experts ( $K$ ) during training. However, this approach transforms the sparse network into a dense one, contradicting the original motivation behind using SMoE for large-scale models.

**Q3: When does VQMoE outperform router methods in terms of robustness? A3: The lower hidden size of FFN.**

Compared to the routing methods, VQMoE achieves competitive performance which only requires 80% number of parameters. Figure 2a

Architecture	FLOPs( $\times 10^{10}$ )	Transformer				Transformer-XL			
		SST-2	SST-5	IMDB	BANKING77	SST-2	SST-5	IMDB	BANKING77
VQMoE	<b>5.6145</b>	<b>82.6</b>	<b>41.1</b>	<b>89.5</b>	<b>84.8</b>	<b>83.3</b>	<b>42.0</b>	<b>89.1</b>	<b>85.3</b>
SMoE	7.7620	82.1	39.5	89.3	82.6	80.8	40.4	88.6	80.2
SMoE/Dropout	7.7620	81.3	39.6	88.9	77.9	81.8	40.0	89.1	77.3
XMoE	7.7620	82.4	39.9	89.0	83.1	81.3	40.3	88.7	82.7
StableMoE	7.7620	82.2	40.4	89.1	82.7	82.5	41.1	88.5	78.6

Table 2: Accuracy of the model after fine-tuned on various datasets. Higher is better, best results are in bold.

and Figure 2b demonstrate the robustness of our method on the Enwik8 and Text8 datasets, respectively.

### 5.3 Parameter-Efficient Fine-Tuning

**Q4: What is the biggest advantage of SMoE, compared to the conventional SMoE? A4: Parameter-Efficient Fine-Tuning.**

We see that the discrete representation that VQMoE learns at the Pretraining stage 5.2 might consist of rich knowledge. To test this hypothesis, we use only the discrete representation for downstream tasks, allowing VQMoE to **save 28%** of computational resources compared to SMoE. Table 2 reports the accuracy of the models fine-tuned on the test sets of various datasets. Overall, we observe that VQMoE demonstrates strong transfer learning capabilities by achieving the highest accuracy on all datasets. Notably, on the more challenging datasets of SST-5 and BANKING77, which have fewer training samples or more classes, we observe larger performance gains from VQMoE versus the remaining baselines (over 5% improvements compared to the second-best method). This result shows that VQMoE can learn a discrete representation that is not only good for pre-training but also exhibits strong transfer capabilities to various downstream tasks.

### 5.4 Vision

**Q5: Can VQMoE compete with SMoE in the Vision domain? A5: Yes.**

To make our performance comparison informative and comprehensive, we consider two kinds of baselines that are fairly comparable to VQMoE: (1) Dense Model (Vision Transformer) (Dosovitskiy et al., 2021); (2) SoftMoE (Puigcerver et al., 2024) - the most advanced MoE in Vision domain. We perform two configurations for training the Mixture of Experts: (1) small - 10 million parameters (10M); (2) large - 110 million parameters (110M). The result at Table 3 shows that VQMoE outperforms both Vision Transformer Dense (Dosovitskiy et al., 2021), SoftMoE (Puigcerver et al., 2024), and other routing methods such as (Dai et al., 2022),

(Chi et al., 2022) on six out of eight tasks across four image classification datasets. We conduct our experiments three times on four datasets (CIFAR-10, CIFAR-100, STL-10, and SVHN) using different seeds, reporting the average results along with the standard deviation. For the large-scale dataset ImageNet-1K, we perform a single run due to resource constraints. The average performance of our method surpasses other baselines and is more stable, as indicated by the low standard deviation.

Architecture # params	Vision Transformer (Small) 10M					Vision Transformer (Large) 110M					Average
	Cifar10	Cifar100	STL-10	SVHN	ImageNet-1K	Cifar10	Cifar100	STL-10	SVHN	ImageNet-1K	
VQMoE	<b>89.7<math>\pm</math>0.4</b>	<b>67.3<math>\pm</math>0.4</b>	<b>66.5<math>\pm</math>0.3</b>	<b>95.6<math>\pm</math>0.1</b>	<b>54.8</b>	<b>92.8<math>\pm</math>0.3</b>	<b>67.0<math>\pm</math>0.5</b>	<b>64.3<math>\pm</math>0.5</b>	<b>96.0<math>\pm</math>0.2</b>	<b>71.3</b>	<b>76.5<math>\pm</math>0.3</b>
SMoE	88.7 $\pm$ 0.2	65.4 $\pm$ 0.5	66.4 $\pm$ 0.1	95.4 $\pm$ 0.1	52.8	85.7 $\pm$ 0.5	55.5 $\pm$ 0.8	64.4 $\pm$ 0.2	94.5 $\pm$ 0.1	71.0	74.0 $\pm$ 1.6
XMoE	88.8 $\pm$ 0.2	65.5 $\pm$ 0.1	66.3 $\pm$ 0.2	95.4 $\pm$ 0.1	52.5	87.1 $\pm$ 0.4	55.9 $\pm$ 0.6	64.6 $\pm$ 0.3	94.1 $\pm$ 0.2	70.8	74.2 $\pm$ 1.1
StableMoE	88.8 $\pm$ 0.1	65.5 $\pm$ 0.1	66.5 $\pm$ 0.2	95.4 $\pm$ 0.1	52.5	84.7 $\pm$ 0.5	55.5 $\pm$ 1.8	64.3 $\pm$ 0.6	94.3 $\pm$ 0.9	70.6	73.8 $\pm$ 1.8
SoftMoE	85.6 $\pm$ 0.3	61.4 $\pm$ 0.2	65.4 $\pm$ 0.2	94.8 $\pm$ 0.1	41.6	80.3 $\pm$ 0.7	42.9 $\pm$ 1.4	63.2 $\pm$ 0.5	93.5 $\pm$ 0.1	68.2	69.7 $\pm$ 1.6
VIT (Dense)	89.0 $\pm$ 0.2	65.7 $\pm$ 0.3	<b>66.6<math>\pm</math>0.2</b>	95.6 $\pm$ 0.1	52.2	92.2 $\pm$ 0.3	60.2 $\pm$ 0.6	64.1 $\pm$ 0.5	96.0 $\pm$ 0.1	71.1	75.3 $\pm$ 0.5

Table 3: Accuracy of models evaluated on vision datasets. Higher is better, the best results are in bold.

### 5.5 In-depth Analysis

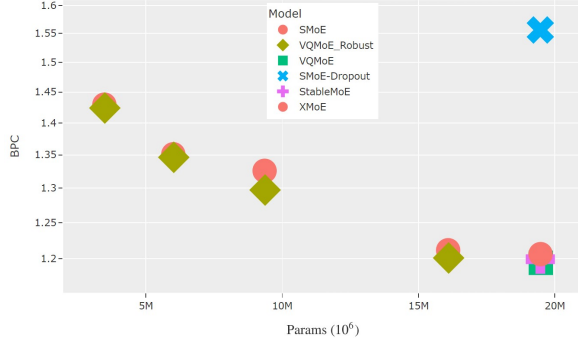
**Consistent Score.** Figure 3a illustrates that expert selections when training SMoE face inconsistent problems. As the Theorem A.3, this inconsistency arises because the router’s coverage rate significantly exceeds that of the Transformer representation. Figure 3a also shows that our method achieves the highest consistency score compared to the SMoE and XMoE models. However, the VQMoE model’s consistency score is around 75%, as our method also requires learning a continuous representation during the Pre-training phase.

**Representation Collapse issue.** To visualize the Representation collapse problem in practice, we apply Principal Component Analysis (PCA) method to reduce from  $d$  dimension of the Transformer to 2D for plotting purposes, thanks to (Chi et al., 2022). Figures 3b and 3c show the expert representations from the pretrained VQMoE and SMoE models. The results suggest that VQMoE experiences less representation collapse in the expert space compared to SMoE. The analysis is in line with the theorem proof at Section 4. However, projecting the  $d$ -dimensional space onto 2D for visualization may lead to information loss.

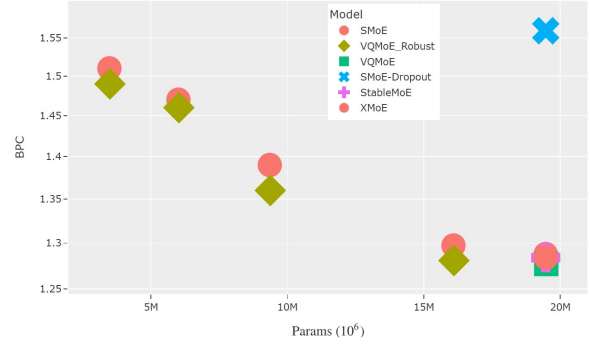
### 5.6 Ablation Study

We examine the effectiveness of VQMoE across various hyper-parameter settings, with all experiments conducted using the base Transformer architecture on the WikiText-103 dataset.

**Vector Quantization Method.** To learn a discrete representation, we research various types of Vector Quantization methods, including VQ-

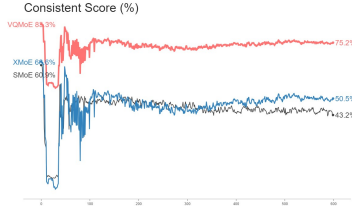


(a) Robust VQMoE Benchmark (Enwik8)

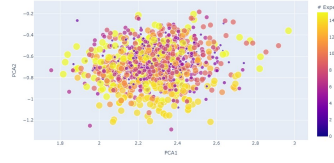


(b) Robust VQMoE Benchmark (Text8)

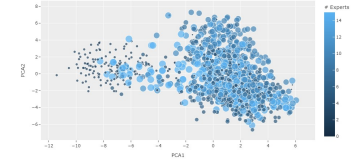
Figure 2: Illustration of the proposed Robust VQMoE architecture for Pre-training on Enwik8 and Text8 dataset. (a) Robust VQMoE architecture achieves the same performance with the routing methods while only using 80% of the parameters on Enwik8 dataset. (b) Robust VQMoE demonstrates robustness on the Text8 dataset. Bits-per-character (BPC) on the Enwik8 and Text8 datasets, and lower is better.



(a) Consistent Score.



(b) VQMoE Representation.



(c) SMOE Representation.

Figure 3: Analysis Inconsistent Expert Selection and Representation Collapse issues when training SMOE. Figure 3a demonstrates consistent score movement from VQMoE, compared with SMOE and XMoE. Figure 3b and Figure 3c visualize the representation by experts in 2D dimension using Principal Component Analysis (PCA) method.

VAE (van den Oord et al., 2017), VQGAN (Yu et al., 2022), LFQ (Yu et al., 2023), and ResidualVQ (Yang et al., 2023). We observe that VQGAN using cosine similarity for distance achieves good and stable results in practice as Figure 5a. Interestingly, VQGAN with lower dimensionality also delivers strong performance and exhibits robustness.

**Number of codebook impact.** The number of codebook entries is a crucial hyperparameter when training Vector Quantization techniques. As shown in Figure 5b, we can see the best performance when the number of codebook entries matches the number of experts. This aligns with the proof by (Dikkala et al., 2023), which demonstrates that in the optimal case, the number of clusters equals the number of experts.

**Sensitiveness of VQ loss contribution  $\alpha$ .** Figure 5c illustrates the impact of  $\alpha$ , which controls the contribution of the Vector Quantization loss to the overall loss. If  $\alpha$  is too high, it leads to a better discrete representation but may negatively affect the final target. Conversely, if  $\alpha$  is too low, it may result in a poor discrete representation. Therefore,

$\alpha$  should be selected based on the data, typically within the range of (0.05, 0.15).

## 6 Conclusion and Future Directions

This study illustrates Vector-Quantized Mixture of Experts (VQMoE), a novel and theoretically-grounded architecture, to overcome challenges in training SMOE such as representation collapse and inconsistency. We evaluate our method on various Pre-training and Fine-tuning tasks, for both language and vision domains. The results show that VQMoE outperforms the routing methods both theoretically and empirically. Furthermore, fine-tuning VQMoE with the discrete representation for downstream tasks could reduce computational resource usage by 28%. We believe that focusing on discrete representation learning will offer a promising strategy for training and testing sparse mixtures of experts (SMoE) at a large scale. Finally, we believe that our approach opens up new research avenues for effectively training SMOE, where cutting-edge techniques in discrete representation learning and vector quantization can be harnessed to enhance their performance.



## Limitations

Our study focuses on enhancing the efficiency and effectiveness of training large language models (LLMs) with SMoE. Although our results are promising, our experiments were restricted to medium-scale datasets and base and large language models due to computational limitations. Consequently, additional empirical evaluations are required to assess the scalability of VQMoE and other SMoE approaches on modern LLMs with up to a few billion parameters.

## Ethics Statement

Despite promising results, training large-scale LLMs remains inherently costly and demands significant computational resources, which must be carefully managed. Additionally, our paper utilized web-sourced data, which is known to contain gender and racial biases, necessitating further efforts to mitigate these negative impacts. Lastly, while our study marks a promising step toward advancing the development of new LLMs, it underscores the need for careful regularization to prevent potential misuse in harmful applications.

## References

Mikel Artetxe, Shruti Bhosale, Naman Goyal, Todor Mihaylov, Myle Ott, Sam Shleifer, Xi Victoria Lin, Jingfei Du, Srinivasan Iyer, Ramakanth Pasunuru, Giri Anantharaman, Xian Li, Shuohui Chen, Halil Akin, Mandeep Baines, Louis Martin, Xing Zhou, Punit Singh Koura, Brian O'Horo, Jeff Wang, Luke Zettlemoyer, Mona Diab, Zornitsa Kozareva, and Ves Stoyanov. 2022. [Efficient large scale language modeling with mixtures of experts](#). *Preprint*, arXiv:2112.10684.

Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. [Efficient intent detection with dual sentence encoders](#). In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 38–45, Online. Association for Computational Linguistics.

Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony Robinson. 2014. [One billion word benchmark for measuring progress in statistical language modeling](#). *Preprint*, arXiv:1312.3005.

Tianlong Chen, Zhenyu Zhang, Ajay Jaiswal, Shiwei Liu, and Zhangyang Wang. 2023a. [Sparse moe as the new dropout: Scaling dense and self-slimmable transformers](#). *Preprint*, arXiv:2303.01610.

Zitian Chen, Yikang Shen, Mingyu Ding, Zhenfang Chen, Hengshuang Zhao, Erik G. Learned-Miller,

and Chuang Gan. 2023b. [Mod-squad: Designing mixtures of experts as modular multi-task learners](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11828–11837.

Zixiang Chen, Yihe Deng, Yue Wu, Quanquan Gu, and Yuanzhi Li. 2022. [Towards understanding the mixture-of-experts layer in deep learning](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 23049–23062. Curran Associates, Inc.

Zewen Chi, Li Dong, Shaohan Huang, Damai Dai, Shuming Ma, Barun Patra, Saksham Singhal, Payal Bajaj, Xia Song, Xian-Ling Mao, Heyan Huang, and Furu Wei. 2022. [On the representation collapse of sparse mixture of experts](#). *Preprint*, arXiv:2204.09179.

Adam Coates, Andrew Ng, and Honglak Lee. 2011. [An Analysis of Single Layer Networks in Unsupervised Feature Learning](#). In *AISTATS*. [https://cs.stanford.edu/~acoates/papers/coatesleeng\\_aistats\\_2011.pdf](https://cs.stanford.edu/~acoates/papers/coatesleeng_aistats_2011.pdf).

Damai Dai, Li Dong, Shuming Ma, Bo Zheng, Zhifang Sui, Baobao Chang, and Furu Wei. 2022. [Stable-moe: Stable routing strategy for mixture of experts](#). *Preprint*, arXiv:2204.08396.

Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc Le, and Ruslan Salakhutdinov. 2019a. [Transformer-XL: Attentive language models beyond a fixed-length context](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988, Florence, Italy. Association for Computational Linguistics.

Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V. Le, and Ruslan Salakhutdinov. 2019b. [Transformer-xl: Attentive language models beyond a fixed-length context](#). *Preprint*, arXiv:1901.02860.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. [Imagenet: A large-scale hierarchical image database](#). In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee.

Nishanth Dikkala, Nikhil Ghosh, Raghu Meka, Rina Panigrahy, Nikhil Vyas, and Xin Wang. 2023. [On the benefits of learning to route in mixture-of-experts models](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9376–9396, Singapore. Association for Computational Linguistics.

Giang Do, Hung Le, and Truyen Tran. 2024. [Simsmoe: Solving representational collapse via similarity measure](#). *Preprint*, arXiv:2406.15883.

Giang Do, Khiem Le, Quang Pham, TrungTin Nguyen, Thanh-Nam Doan, Bint T. Nguyen, Chenghao Liu, Savitha Ramasamy, Xiaoli Li, and Steven Hoi. 2023. [Hyperrouter: Towards efficient training and](#)

714	<a href="#">inference of sparse mixture of experts.</a>	<i>Preprint</i> ,	770
715	<a href="#">arXiv:2312.07035.</a>		771
716	Alexey Dosovitskiy, Lucas Beyer, Alexander		772
717	Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,		773
718	Thomas Unterthiner, Mostafa Dehghani, Matthias		774
719	Minderer, Georg Heigold, Sylvain Gelly, Jakob		775
720	Uszkoreit, and Neil Houlsby. 2021. <a href="#">An image</a>		776
721	<a href="#">is worth 16x16 words: Transformers for image</a>		777
722	<a href="#">recognition at scale.</a>	<i>Preprint</i> , arXiv:2010.11929.	778
723	Nan Du, Yanping Huang, Andrew M. Dai, Simon Tong,		779
724	Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun,		780
725	Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret		781
726	Zoph, Liam Fedus, Maarten Bosma, Zongwei Zhou,		782
727	Tao Wang, Yu Emma Wang, Kellie Webster, Marie		783
728	Pellat, Kevin Robinson, Kathleen Meier-Hellstern,		784
729	Toju Duke, Lucas Dixon, Kun Zhang, Quoc V Le,		785
730	Yonghui Wu, Zhifeng Chen, and Claire Cui. 2022.		786
731	<a href="#">Glam: Efficient scaling of language models with</a>		787
732	<a href="#">mixture-of-experts.</a>	<i>Preprint</i> , arXiv:2112.06905.	788
733	William Fedus, Barret Zoph, and Noam Shazeer. 2022.		789
734	<a href="#">Switch transformers: Scaling to trillion parameter</a>		790
735	<a href="#">models with simple and efficient sparsity.</a>	<i>Preprint</i> ,	791
736	<a href="#">arXiv:2101.03961.</a>		792
737	Mor Geva, Roei Schuster, Jonathan Berant, and Omer		793
738	Levy. 2021. <a href="#">Transformer feed-forward layers are key-</a>		794
739	<a href="#">value memories.</a>	In <i>Proceedings of the 2021 Confer-</i>	795
740	<i>ence on Empirical Methods in Natural Language Pro-</i>		796
741	<i>cessing</i> , pages 5484–5495, Online and Punta Cana,		797
742	Dominican Republic. Association for Computational		798
743	Linguistics.		799
744	Yunhao Gou, Zhili Liu, Kai Chen, Lanqing Hong, Hang		800
745	Xu, Aoxue Li, Dit-Yan Yeung, James T. Kwok, and		801
746	Yu Zhang. 2024. <a href="#">Mixture of cluster-conditional</a>		802
747	<a href="#">lora experts for vision-language instruction tuning.</a>		803
748	<i>Preprint</i> , arXiv:2312.12379.		804
749	David Ha, Andrew Dai, and Quoc V. Le. 2016. <a href="#">Hyper-</a>		805
750	<a href="#">networks.</a>	<i>Preprint</i> , arXiv:1609.09106.	806
751	Weihua Hu, Takeru Miyato, Seiya Tokui, Eiichi Mat-		807
752	sumoto, and Masashi Sugiyama. 2017. <a href="#">Learning</a>		808
753	<a href="#">discrete representations via information maximizing</a>		809
754	<a href="#">self-augmented training.</a>	In <i>Proceedings of the 34th</i>	810
755	<i>International Conference on Machine Learning</i> , vol-		811
756	ume 70 of <i>Proceedings of Machine Learning Re-</i>		812
757	<i>search</i> , pages 1558–1567. PMLR.		813
758	Changho Hwang, Wei Cui, Yifan Xiong, Ziyue Yang,		814
759	Ze Liu, Han Hu, Zilong Wang, Rafael Salas, Jithin		815
760	Jose, Prabhat Ram, Joe Chau, Peng Cheng, Fan		816
761	Yang, Mao Yang, and Yongqiang Xiong. 2023. <a href="#">Tu-</a>		817
762	<a href="#">tel: Adaptive mixture-of-experts at scale.</a>	<i>Preprint</i> ,	818
763	<a href="#">arXiv:2206.03382.</a>		819
764	Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan,		820
765	and Geoffrey E. Hinton. 1991a. <a href="#">Adaptive mixtures</a>		821
766	<a href="#">of local experts.</a>	<i>Neural Computation</i> , 3(1):79–87.	822
767	Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan,		823
768	and Geoffrey E. Hinton. 1991b. <a href="#">Adaptive mixtures</a>		824
769	<a href="#">of local experts.</a>	<i>Neural Computation</i> , 3(1):79–87.	
	Albert Q. Jiang, Alexandre Sablayrolles, Antoine		
	Roux, Arthur Mensch, Blanche Savary, Chris		
	Bamford, Devendra Singh Chaplot, Diego de las		
	Casas, Emma Bou Hanna, Florian Bressand, Gi-		
	anna Lengyel, Guillaume Bour, Guillaume Lam-		
	ple, L��lio Renard Lavaud, Lucile Saulnier, Marie-		
	Anne Lachaux, Pierre Stock, Sandeep Subramanian,		
	Sophia Yang, Szymon Antoniak, Teven Le Scao,		
	Th��ophile Gerv��t, Thibaut Lavril, Thomas Wang,		
	Timoth��e Lacroix, and William El Sayed. 2024. <a href="#">Mix-</a>		
	<a href="#">tral of experts.</a>	<i>Preprint</i> , arXiv:2401.04088.	
	Michael Jordan and Robert Jacobs. 1994. Hierarchical		
	mixtures of experts and the. <i>Neural computation</i> ,		
	6:181–.		
	Jean Kaddour, Joshua Harris, Maximilian Mozes, Her-		
	bie Bradley, Roberta Raileanu, and Robert McHardy.		
	2023. <a href="#">Challenges and applications of large language</a>		
	<a href="#">models.</a>	<i>Preprint</i> , arXiv:2307.10169.	
	Diederik P. Kingma and Jimmy Ba. 2017. <a href="#">Adam:</a>		
	<a href="#">A method for stochastic optimization.</a>	<i>Preprint</i> ,	
	<a href="#">arXiv:1412.6980.</a>		
	Diederik P Kingma and Max Welling. 2022. <a href="#">Auto-</a>		
	<a href="#">encoding variational bayes.</a>	<i>Preprint</i> ,	
	<a href="#">arXiv:1312.6114.</a>		
	Jakub Krajewski, Jan Ludziejewski, Kamil Adam-		
	czewski, Maciej Pi��ro, Micha�� Krutul, Szymon		
	Antoniak, Kamil Ciebiera, Krystian Kr��l, Tomasz		
	Odrzyg��������, Piotr Sankowski, Marek Cygan, and Se-		
	bastian Jaszczur. 2024. <a href="#">Scaling laws for fine-grained</a>		
	<a href="#">mixture of experts.</a>	<i>Preprint</i> , arXiv:2402.07871.	
	Alex Krizhevsky. 2009. Learning multiple layers of		
	features from tiny images. Technical report, UoT.		
	Yoohwan Kwon and Soo-Whan Chung. 2023. <a href="#">Mole :</a>		
	<a href="#">Mixture of language experts for multi-lingual auto-</a>		
	<a href="#">matic speech recognition.</a>	In <i>ICASSP 2023 - 2023</i>	
	<i>IEEE International Conference on Acoustics, Speech</i>		
	<i>and Signal Processing (ICASSP)</i> , pages 1–5.		
	Margaret Li, Suchin Gururangan, Tim Dettmers, Mike		
	Lewis, Tim Althoff, Noah A. Smith, and Luke Zettle-		
	moyer. 2022. <a href="#">Branch-train-merge: Embarrassingly</a>		
	<a href="#">parallel training of expert language models.</a>	<i>Preprint</i> ,	
	<a href="#">arXiv:2208.03306.</a>		
	Bin Lin, Zhenyu Tang, Yang Ye, Jiaxi Cui, Bin Zhu,		
	Peng Jin, Jinfa Huang, Junwu Zhang, Munan Ning,		
	and Li Yuan. 2024. <a href="#">Moe-llava: Mixture of ex-</a>		
	<a href="#">perts for large vision-language models.</a>	<i>Preprint</i> ,	
	<a href="#">arXiv:2401.15947.</a>		
	Ilya Loshchilov and Frank Hutter. 2017. <a href="#">Sgdr: Stochas-</a>		
	<a href="#">tic gradient descent with warm restarts.</a>	<i>Preprint</i> ,	
	<a href="#">arXiv:1608.03983.</a>		
	Andrew L. Maas, Raymond E. Daly, Peter T. Pham,		
	Dan Huang, Andrew Y. Ng, and Christopher Potts.		
	2011. <a href="#">Learning Word Vectors for Sentiment Analy-</a>		
	<a href="#">sis.</a>	In <i>Proceedings of the 49th Annual Meeting of the</i>	
	<i>Association for Computational Linguistics: Human</i>		

825	<i>Language Technologies</i> , pages 142–150, Portland,	878
826	Oregon, USA. Association for Computational Lin-	879
827	guistics.	880
828	Matt Mahoney. 2011. <a href="#">Large text compression bench-</a>	881
829	<a href="#">mark</a> .	882
830	Chengzhi Mao, Lu Jiang, Mostafa Dehghani, Carl Von-	883
831	drick, Rahul Sukthankar, and Irfan Essa. 2022. <a href="#">Dis-</a>	884
832	<a href="#">crete representations strengthen vision transformer</a>	
833	<a href="#">robustness</a> . <i>Preprint</i> , arXiv:2111.10493.	
834	Fabian Mentzer, David Minnen, Eirikur Agustsson, and	
835	Michael Tschanen. 2023. <a href="#">Finite scalar quantization:</a>	
836	<a href="#">Vq-vae made simple</a> . <i>Preprint</i> , arXiv:2309.15505.	
837	Stephen Merity, Caiming Xiong, James Bradbury, and	
838	Richard Socher. 2016. <a href="#">Pointer sentinel mixture mod-</a>	
839	<a href="#">els</a> . <i>Preprint</i> , arXiv:1609.07843.	
840	Yuval Netzer, Tao Wang, Adam Coates, Alessandro	
841	Bissacco, Bo Wu, and Andrew Y Ng. 2011. Reading	
842	digits in natural images with unsupervised feature	
843	learning. <i>NIPS Workshop</i> .	
844	Quang Pham, Giang Do, Huy Nguyen, TrungTin	
845	Nguyen, Chenghao Liu, Mina Sartipi, Binh T.	
846	Nguyen, Savitha Ramasamy, Xiaoli Li, Steven Hoi,	
847	and Nhat Ho. 2024. <a href="#">Competesmoe – effective train-</a>	
848	<a href="#">ing of sparse mixture of experts via competition</a> .	
849	<i>Preprint</i> , arXiv:2402.02526.	
850	Joan Puigcerver, Carlos Riquelme, Basil Mustafa, and	
851	Neil Houlsby. 2024. <a href="#">From sparse to soft mixtures of</a>	
852	<a href="#">experts</a> . <i>Preprint</i> , arXiv:2308.00951.	
853	Carlos Riquelme, Joan Puigcerver, Basil Mustafa,	
854	Maxim Neumann, Rodolphe Jenatton, André Su-	
855	sano Pinto, Daniel Keysers, and Neil Houlsby.	
856	2021a. <a href="#">Scaling vision with sparse mixture of experts</a> .	
857	<i>Preprint</i> , arXiv:2106.05974.	
858	Carlos Riquelme, Joan Puigcerver, Basil Mustafa,	
859	Maxim Neumann, Rodolphe Jenatton, André Su-	
860	sano Pinto, Daniel Keysers, and Neil Houlsby. 2021b.	
861	<a href="#">Scaling vision with sparse mixture of experts</a> . In	
862	<i>Advances in Neural Information Processing Systems</i> ,	
863	volume 34, pages 8583–8595. Curran Associates,	
864	Inc.	
865	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarczy,	
866	Andy Davis, Quoc Le, Geoffrey Hinton, and	
867	Jeff Dean. 2017. <a href="#">Outrageously large neural net-</a>	
868	<a href="#">works: The sparsely-gated mixture-of-experts layer</a> .	
869	<i>Preprint</i> , arXiv:1701.06538.	
870	Sheng Shen, Le Hou, Yanqi Zhou, Nan Du, Shayne	
871	Longpre, Jason Wei, Hyung Won Chung, Barret	
872	Zoph, William Fedus, Xinyun Chen, Tu Vu, Yuxin	
873	Wu, Wuyang Chen, Albert Webson, Yunxuan Li, Vin-	
874	cent Zhao, Hongkun Yu, Kurt Keutzer, Trevor Darrell,	
875	and Denny Zhou. 2023a. <a href="#">Mixture-of-experts meets</a>	
876	<a href="#">instruction tuning: a winning combination for large</a>	
877	<a href="#">language models</a> . <i>Preprint</i> , arXiv:2305.14705.	
	Sheng Shen, Zhewei Yao, Chunyuan Li, Trevor Dar-	878
	rell, Kurt Keutzer, and Yuxiong He. 2023b. <a href="#">Scal-</a>	879
	<a href="#">ing vision-language models with sparse mixture of</a>	880
	<a href="#">experts</a> . In <i>Findings of the Association for Computa-</i>	881
	<i>tional Linguistics: EMNLP 2023</i> , pages 11329–	882
	11344, Singapore. Association for Computational	883
	Linguistics.	884
	Richard Socher, Alex Perelygin, Jean Wu, Jason	885
	Chuang, Christopher D. Manning, Andrew Ng, and	886
	Christopher Potts. 2013. <a href="#">Recursive Deep Models for</a>	887
	<a href="#">Semantic Compositionality Over a Sentiment Tree-</a>	888
	<a href="#">bank</a> . In <i>Proceedings of the 2013 Conference on</i>	889
	<i>Empirical Methods in Natural Language Processing</i> ,	890
	pages 1631–1642, Seattle, Washington, USA. Asso-	891
	ciation for Computational Linguistics.	892
	Robin Strudel, Ricardo Garcia, Ivan Laptev, and	893
	Cordelia Schmid. 2021. <a href="#">Segmenter: Transformer for</a>	894
	<a href="#">semantic segmentation</a> . <i>Preprint</i> , arXiv:2105.05633.	895
	Aaron van den Oord, Oriol Vinyals, and koray	896
	kavukcuoglu. 2017. <a href="#">Neural discrete representation</a>	897
	<a href="#">learning</a> . In <i>Advances in Neural Information Pro-</i>	898
	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	899
	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	900
	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	901
	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all</a>	902
	<a href="#">you need</a> . In <i>Advances in Neural Information Pro-</i>	903
	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.	904
	Wenxuan Wang, Guodong Ma, Yuke Li, and Binbin	905
	Du. 2023. <a href="#">Language-routing mixture of experts for</a>	906
	<a href="#">multilingual and code-switching speech recognition</a> .	907
	<i>Preprint</i> , arXiv:2307.05956.	908
	Fuzhao Xue, Zian Zheng, Yao Fu, Jinjie Ni, Zang-	909
	wei Zheng, Wangchunshu Zhou, and Yang You.	910
	2024. <a href="#">Openmoe: An early effort on open</a>	911
	<a href="#">mixture-of-experts language models</a> . <i>Preprint</i> ,	912
	arXiv:2402.01739.	913
	Dongchao Yang, Songxiang Liu, Rongjie Huang,	914
	Jinchuan Tian, Chao Weng, and Yuxian Zou.	915
	2023. <a href="#">Hifi-codec: Group-residual vector quan-</a>	916
	<a href="#">tization for high fidelity audio codec</a> . <i>Preprint</i> ,	917
	arXiv:2305.02765.	918
	Hanrong Ye and Dan Xu. 2023. Taskexpert: Dynam-	919
	ically assembling multi-task representations with	920
	memorial mixture-of-experts. In <i>Proceedings of the</i>	921
	<i>IEEE/CVF International Conference on Computer</i>	922
	<i>Vision (ICCV)</i> , pages 21828–21837.	923
	Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruom-	924
	ing Pang, James Qin, Alexander Ku, Yuanzhong Xu,	925
	Jason Baldridge, and Yonghui Wu. 2022. <a href="#">Vector-</a>	926
	<a href="#">quantized image modeling with improved vqgan</a> .	927
	<i>Preprint</i> , arXiv:2110.04627.	928
	Lijun Yu, Yong Cheng, Kihyuk Sohn, José Lezama,	929
	Han Zhang, Huiwen Chang, Alexander G. Haupt-	930
	mann, Ming-Hsuan Yang, Yuan Hao, Irfan Essa, and	931
	Lu Jiang. 2023. <a href="#">Magvit: Masked generative video</a>	932
	<a href="#">transformer</a> . <i>Preprint</i> , arXiv:2212.05199.	933



Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. 2018. [Semantic understanding of scenes through the ade20k dataset](#). *Preprint*, arXiv:1608.05442.

Yanqi Zhou, Nan Du, Yanping Huang, Daiyi Peng, Chang Lan, Da Huang, Siamak Shakeri, David So, Andrew Dai, Yifeng Lu, Zhifeng Chen, Quoc Le, Claire Cui, James Laudon, and Jeff Dean. 2024. [Brainformers: Trading simplicity for efficiency](#). *Preprint*, arXiv:2306.00008.

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-of-experts with expert choice routing](#). In *Advances in Neural Information Processing Systems*, volume 35, pages 7103–7114. Curran Associates, Inc.

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](#). *Preprint*, arXiv:2202.08906.

## A Appendix

### Supplementary Material for “On the effectiveness of discrete representations in sparse mixture of experts”

This document is organized as follows. Appendix A.1 provides a detailed theoretical analysis of the SMOE Router. Appendix A.2 presents additional experimental results demonstrating the effectiveness of our method compared to the baselines. Finally, Appendix A.3 offers an in-depth analysis of representation collapse, while Appendix A.4 details the implementation aspects.

#### A.1 Theory Analysis for SMOE Router

##### A.1.1 Optimal Experts Selection

**Problem settings.** We consider an MoE layer with each expert being an MLP layer which is trained by gradient descent and input data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  generated from a data distribution  $\mathcal{D}$ . Same as (Chen et al., 2022); (Dikkala et al., 2023), we assume that the MoE input exhibits cluster properties, meaning the data is generated from  $K$  distinct clusters  $(C_1, C_2, \dots, C_k)$ .

**Definition A.1 (Consistent Router)** A sequence of points  $x_1, x_2, \dots, x_n$  and a corresponding sequence of clusters  $C_1, C_2, \dots, C_k$  are said to be **consistent** if, for every point  $x_p \in C_i$ , the condition

$$\text{dist}(x_p, u_i) \leq \min_{j \neq i} \text{dist}(x_p, u_j)$$

is satisfied, where  $\text{dist}(a, b)$  denotes the distance between  $a$  and  $b$ , and  $u_i$  is the center of cluster  $C_i$ .

**Definition A.2 (Inconsistent Router)** A sequence of points  $x_1, x_2, \dots, x_n$  and a corresponding sequence of clusters  $C_1, C_2, \dots, C_k$  are said to be **inconsistent** if there exists a point  $x_p \in C_i$  such that

$$\text{dist}(x_p, u_i) > \min_{j \neq i} \text{dist}(x_p, u_j),$$

where  $\text{dist}(a, b)$  represents the distance between  $a$  and  $b$ , and  $u_i$  is the center of cluster  $C_i$ .

Inspired by (Dikkala et al., 2023), we conceptualize the router in Sparse Mixture of Experts as a clustering problem. This leads us to define a consistent router in Definition A.1. Furthermore, we introduce a definition for an inconsistent router in SMOE as outlined in Definition A.2, along with the concept of inconsistent expert selection presented in Theorem A.3 during the training of SMOE.

##### Theorem A.3 (Inconsistent Experts Selection)

Let  $f_{MHA}$  be a multi-head attention (MHA) function producing an output  $x \in \mathbb{R}^{n \times d}$ , and consider  $N$  experts with embeddings  $e_i$  for expert  $i$  where  $i \in [1, N]$ . Assume that  $f_{MHA}$  converges at step  $t_m$ , while the expert embeddings  $e$  converge at step  $t_e$ , with  $t_m \gg t_e$ . For each output  $x$ , the expert  $K \in [1, N]$  is selected such that

$$K = \arg \min_{j \in [1, N]} \text{dist}(x, e_j).$$

Under these conditions, the expert embeddings  $e$  form an inconsistent routing mechanism.

The proof of Theorem A.3 is given in Appendix A, and we have the following insights. Theorem A.3 implies that an expert selection process by a router as the conventional SMOE leads to the inconsistent router. Indeed, the router layer is designed as a simple linear layer,  $x$  is the output of MHA function in practice. In practice, an SMOE router is significantly simpler than the MHA function. Consequently, this design leads to the router functioning as an inconsistent router, contributing to the representation collapse issue and instability during training.

##### Proposition A.4 (Optimal Experts Selection)

Given input data partitioned into  $k$  clusters  $(C_1, C_2, \dots, C_k)$  and a mixture of experts (MoE) layer with  $k$  experts  $(E_1, E_2, \dots, E_k)$ , the assignment of each cluster  $C_i$  to expert  $E_i$  for  $i \in [1, k]$  constitutes an optimal expert selection solution.



Proposition A.4 demonstrates that if we are given a clustering structure as input, assigning each part of the input to its corresponding expert results in an optimal expert selection. This implies that learning a discrete representation and directing each component to the appropriate expert yields an optimal solution. The proof of Proposition A.4 can be found in Appendix A.

### A.1.2 Proof of Theorem A.3

In this proof, we use contradiction to establish the theorem. Assume that the expert embeddings  $e$  form a consistent router. By Definition A.1, we have:

$$\text{dist}(x_p, u_i) \leq \min(\text{dist}(x_p, C_j)),$$

where  $u_i$  is the representation corresponding to the closest expert  $e_i$ .

According to (Chi et al., 2022), projecting information from a hidden representation space  $\mathcal{R}^d$  to the expert dimension  $N$  leads to representation collapse. Now, consider three experts  $x, y, z$  whose embeddings  $e_x, e_y, e_z$  collapse. Without loss of generality, assume that  $e_y$  lies between  $e_x$  and  $e_z$  in the embedding space. Then, we have:

$$\begin{aligned} \text{dist}(y, u_y) &\leq \min(\text{dist}(x, e_x), \text{dist}(y, e_y), \text{dist}(z, e_z)) \\ &\leq \text{dist}(e_x, e_z). \end{aligned} \quad (11)$$

Let  $t_e$  denote the step at which the embeddings  $e_x$  and  $e_z$  converge, and  $t_m$  denote the step at which the Multi-Head Attention (MHA) module converges. From step  $t_e$ , it follows that:

$$\lim_{t_e \rightarrow t_m} \text{dist}(y, u_y) = \lim_{t_e \rightarrow t_m} \text{dist}(e_x, e_z) = 0.$$

Thus,  $y$  (the output of MHA) converges at step  $t_e$ .

This directly contradicts the assumption that the MHA converges at step  $t_m$ , where  $t_e \ll t_m$ .

### A.1.3 Proof of Proposition A.4

We use contradiction to prove the proposition. Assume that, at training step  $t$ , there exists a set of pairs  $(C_i, E_j)$  such that  $i \neq j$ . Let  $x_1, x_2, \dots, x_k$  represent a sequence of inputs sampled from  $K$  clusters. From step  $t_0$  to step  $t_{k-1}$ , each pair  $(x_j, E_j)$ , where  $j \in [1, k]$ , is updated using the following gradient descent equation:

$$W_{E_j}^{l+1} = W_{E_j}^l - \eta \mathcal{J}(x_j),$$

where  $W_{E_j}^l$  is the weight of expert  $E_j$  at iteration  $l$ ,  $\mathcal{J}(x_j)$  is the Jacobian matrix with respect to input  $x_j$ , and  $\eta$  is the learning rate.

Let  $\mathcal{L}$  denote the loss function during the training process described by Equation 6. After  $t_k$  training steps, the following condition holds:

$$E_j(x_j) = \min_{c \in [1, k]} E_j(x_c).$$

Under the assumption of contradiction, there exists a set of pairs

$$\sum_{i,j=1; i \neq j}^K (C_i, E_j)$$

where the loss function  $\mathcal{L}$  is minimized. However, by definition of the loss minimization process, the inequality

$$\sum_{i=1}^K (C_i, E_i) \leq \sum_{i,j=1; i \neq j}^K (C_i, E_j)$$

must hold.

This leads to a contradiction with our initial assumption.

## A.2 Additional Experiment Results

**Q6: Can VQMoE learn Discrete Representation Only from scratch? A6: Yes for small and medium scale, but no for large scale.**

The answer is yes for small and medium-scale models. However, training a discrete representation-only approach is feasible primarily for small to medium-scale models with a moderately sized dataset. The results of the *Transformer-XL* model in Table 4 on the Enwik8 dataset support this observation. As the model scales up, relying solely on discrete representation reaches its limitations, leading to performance below the SMOE baselines.

**Q7: Can VQMoE outperform the clustering-based approach such as KMean? A7: Yes.**

We explored a clustering-based approach similar to MoCLE (Gou et al., 2024) but found it unsuitable for our method. Unlike MoCLE, Vector Quantization allows the model greater flexibility in learning cluster representations during training, making it more competitive in practical applications. The training results using the *Transformer-XL* model on the Enwik8 dataset are presented in Table 5.

**Q8: Can VQMoE contribute to AI real-world applications? A8: Yes.**

Scale	TopK	# Experts	SMoE	VQMoE (Discrete Only)
Base 20M-50K Steps	1	16	1.28	<b>1.25</b>
	2	16	1.26	-
	4	16	1.26	-
	8	16	1.27	-
	16	16	1.27	-
Base 20M-100K Steps	1	16	1.22	<b>1.18</b>
	2	16	1.20	-
	4	16	1.21	-
	8	16	1.21	-
	16	16	1.21	-
Large (210M)	1	64	<b>1.12</b>	1.14
	2	64	1.09	-
	4	64	1.09	-
	8	64	1.09	-
	16	64	1.10	-
	32	64	1.10	-
	64	64	1.12	-

Table 4: Performance comparison of SMoE and VQMoE (Discrete Only) on the *Enwik8* (BPC) dataset.

Scale	TopK	# Experts	SMoE	MoCLE	VQMoE
Base 20M-50K Steps	1	16	1.28	1.29	<b>1.25</b>
	2	16	1.26	1.28	-
	4	16	1.26	1.28	-
	8	16	1.27	1.28	-
	16	16	1.27	1.28	-

Table 5: Performance comparison of VQMoE and MoCLE (Clustering approach) on the *Enwik8* (BPC) dataset.

We found that VQMoE can directly benefit real-world AI applications, such as image segmentation, demonstrating its strong generalization capabilities. Specifically, our method outperforms both the baseline and dense models in terms of Mean Accuracy and mIoU metrics on the ADE20K dataset (Zhou et al., 2018) using the Segmenter model (Strudel et al., 2021). Detailed results are provided in Table 6.

### A.3 Representation Collapse Analysis

To illustrate Theorem A.3, we perform a language model task as described in Section A.4.2, examining the movement of Expert Input Representation in Figure 4a and Expert Embedding (router) in Figure 4b. We analyze the dynamics of the expert input representations by tracking their changes across training iterations. The results indicate that the inputs to the experts become increasingly divergent over time. This divergence suggests that the model learns to represent the data in a more specialized and diverse manner, allowing each expert to focus on distinct features or patterns within the data. Similarly, we track the changes in expert embeddings (router) throughout the training process. However, the trend is the opposite: the expert embeddings appear to converge quickly, stabilizing

Model	ViT	SoftMoe	SMoE	StableMoE	XMoe	VQMoE	Metrics
Segmenter	20.8 15.0	19.0 14.0	23.1 15.5	22.4 16.0	22.3 15.7	<b>23.4</b> <b>16.6</b>	Mean accuracy mIoU

Table 6: Comparison of VQMoE versus the baselines on the ADE20K dataset.

around 10,000 iterations. The findings align with our assumption stated in Theorem A.3, indicating that Expert Embedding converges more quickly than Expert Input Representation. These results provide further evidence supporting the Theorem A.3.

### A.4 Experiments implementation details

This section provides detailed parameters of our experiments in Section 5.

#### A.4.1 General Settings

The experiments are based on the publicly available SMoE-Dropout implementation (Chen et al., 2023a)<sup>1</sup>. However, the pre-training was conducted on two H100 GPUs, so results might differ when using parallel training on multiple GPUs.

#### A.4.2 Pre-training Experiments

Table 7 provides the detailed configurations for pre-training Transformer (Vaswani et al., 2017), Transformer-XL (Dai et al., 2019b) on Enwik8, Text8, WikiText-103, and One Billion Word.

Dataset	Input length	Batch size	Optimizer	Lr	# Training Step
Enwik8	512	48	Adam	3.5e-4	100k
Text	512	48	Adam	3.5e-4	100k
WikiText-103	512	22	Adam	3.5e-4	100k
One Billion Word	512	11	Adam	3.5e-4	100k

Table 7: Hyperparameter settings for pre-training experiments on Enwik8, Text8, WikiText-103, and One Billion Word.

Dataset	Input length	Batch size	Optimizer	Lr	# Epochs
SST-2	512	16	Adam	1e-4	5
SST-5	512	16	Adam	1e-4	5
IMDB	512	4	Adam	1e-4	5
BANKING77	512	16	Adam	1e-4	5

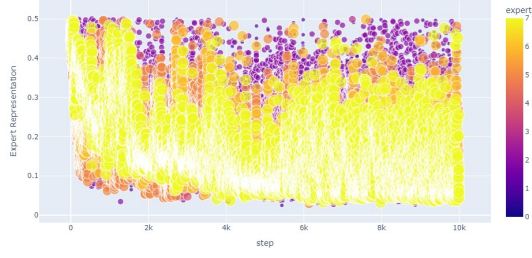
Table 8: Detail settings for fine-tuning experiments on the evaluation datasets.

#### A.4.3 Fine-tuning Experiments

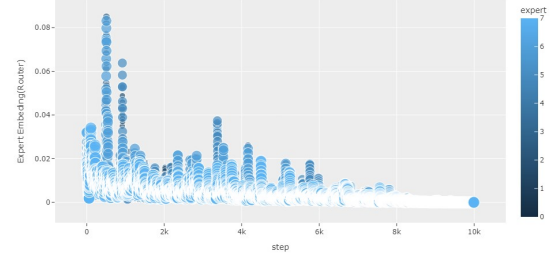
For fine-tuning experiments, we employ the identical model architecture as in pre-training. Table

<sup>1</sup><https://github.com/VITA-Group/Random-MoE-as-Dropout>

8 presents the detailed configurations utilized for fine-tuning experiments on SST-2, SST-5, IMDB, and BANKING77 datasets. We start with the pre-trained checkpoint of the base model on enwik8, remove the final layer, and replace it with two randomly initialized fully connected layers to serve as the classifier for each fine-tuning dataset. All methods are fine-tuned for 5,000 steps with a uniform learning rate.

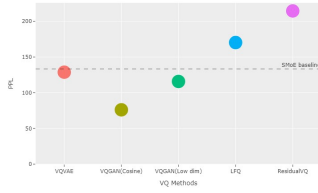


(a) Training Input Token Representations.

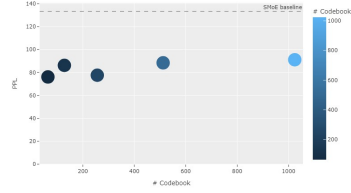


(b) Training Router Representation (Expert embedding).

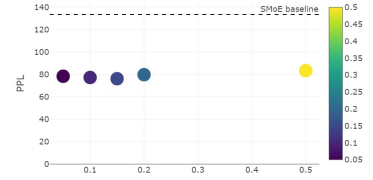
Figure 4: Comparison of Token Representation and Expert Representation across Training Iteration.



(a) Vector Quantization method.



(b) Number of codebook.



(c) Impact of  $\alpha$  for VQMoE.

Figure 5: Pre-training small Transformer-XL on WikiText-103 across different hyperparameters.