# Multi-Head Mixture-of-Experts

Xun Wu [1 2]  Shaohan Huang [2]  Wenhui Wang [2]  Furu Wei [2]

## Abstract

Sparse Mixtures of Experts (SMoE) scales model capacity without significant increases in training and inference costs. However, it exhibits two issues: (1) *Low expert activation*, where only a small subset of experts are activated for optimization, leading to suboptimal performance and limiting its effectiveness in learning a larger number of experts in complex tasks. (2) *Lack of fine-grained analytical capabilities* for multiple semantic concepts within individual tokens. In this paper, we propose Multi-Head Mixture-of-Experts (MH-MoE). MH-MoE employs a multi-head mechanism to split each input token into multiple sub-tokens. Then these sub-tokens are assigned to and processed by a diverse set of experts in parallel, and seamlessly reintegrated into the original token form. The above operations enables MH-MoE to collectively attend to information from various representation spaces within different experts to deepen context understanding while significantly enhancing expert activation. It's worth noting that our MH-MoE is straightforward to implement and decouples from other SMoE frameworks, making it easy to integrate with these frameworks for enhanced performance. Extensive experimental results across three tasks: English-focused language modeling, Multi-lingual language modeling and Masked multi-modality modeling tasks, demonstrate the effectiveness of MH-MoE. Our code are available at https://github.com/yushuiwx/MH-MoE.

## 1. Introduction

Large capacity models, such as Large Language Models (LLMs) (Zhao et al., 2023; Pham et al., 2023; Chung et al., 2022; OpenAI, 2023) and Large Multi-modal Mod-
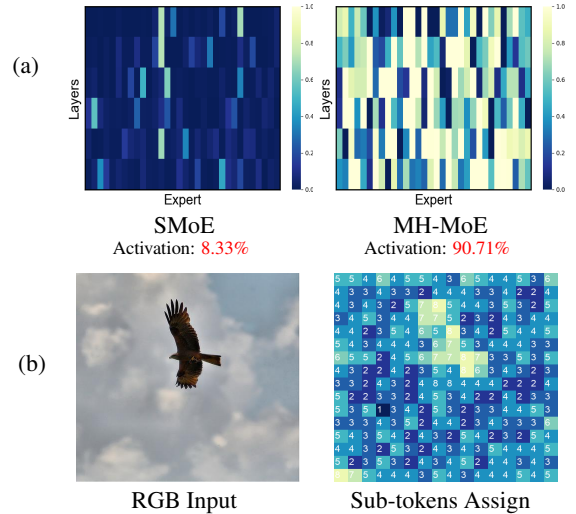
*Figure 1.* (a) **Expert activation distribution** on XNLI (Conneau et al., 2018), encompassing 6 parallel expert layers with 32 experts per layer. SMoE has many "dead" experts (dark) which are not activated, while MH-MoE leading to significantly increased usage of these experts. Experts activation ratio is determined by calculating the ratio of each expert's selection frequency in each MoE layer to the total number of tokens, where those exceeding a threshold (<1) are considered activated. (b) **MH-MoE showcases finer-grained understanding** by distributing sub-tokens split from semantically-rich patches to more distinct experts to capture semantic information. Brighter regions indicate that sub-tokens from this patch are distributed to a greater number of diverse experts, while darker regions indicate that sub-tokens are assigned to more of the same experts.

els (LMMs) (Wang et al., 2022; Peng et al., 2023), have demonstrated their efficacy across various domains and tasks. To further enhance performance, a reliable approach involves scaling up these models by augmenting the parameter count (Fedus et al., 2022). But for most of these densely-activated large-capacity models (referred to as Dense models), which utilize all their parameters to process all inputs, the extremely large size of these models significantly reduces inference speed, further limiting their practicality.

A promising alternative, facilitating model scalability while mitigating the burdensome computational costs, resides in Sparse Mixtures of Experts (SMoE) (Shazeer et al., 2017b; Du et al., 2021; Chi et al., 2022; Clark et al., 2022).
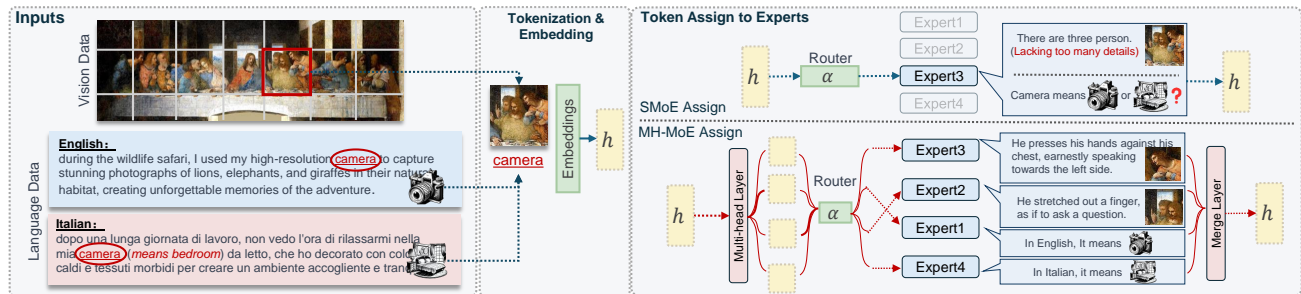
*Figure 2.* **Workflow for MH-MoE on both vision and language data**. For vision data, different heads routed to different experts try to capture different aspects of details within patches and relations between patches. For language data, different heads attend to capture the varying contexts of false cognates across different languages (e.g., Italian and English) or polysemous words within a single language.

In contrast to Dense model, SMoE contains parallel feed-forward neural networks (referred to as experts) within each building block, and strategically activates distinct experts for specific input tokens via a router, thereby yielding noteworthy efficiency enhancements. GShard (Lepikhin et al., 2020) scales a Dense model from 2B to 600B parameters with lower training costs than a 100B Dense model. And recently, Mixtral 8×7B (Jiang et al., 2024), a SMoE model containing 8 experts (7B parameter in total) is shown to outperform or matches LLaMA-2 70B (Touvron et al., 2023) and GPT-3.5.

Despite its success, SMoE has some drawbacks: (1) *Low experts activation*, which means that only a small subset of experts are activated during optimization and inference, e.g., **8.33%** activation ratio shown in Figure 1 (a), while the majority of them are not used at all (see the dark area). As a result, SMoE fails to utilize the full expressive power of these experts, especially when the number of experts is large, which significantly limits effectiveness and scalability of SMoE. (2) *Absence of fine-grained analytical capabilities.* The current tokenization patterns impose limitations on the model's capacity to grasp multiple semantic interpretations linked to individual tokens. In the context of visual data, dividing images into patches for tokenization may either neglect finer image details when using larger patches or escalate computational requirements when employing smaller ones. For language data, the tokenization of false cognates across different languages or polysemous words within a single language results in them being represented by the same tokens, despite carrying distinct meanings. This can subsequently lead to confusion within the models.

To tackle the above issues, we propose Multi-Head Mixture-of-Experts (MH-MoE). The workflow of MH-MoE is illustrated in Figure 2. By employing a multi-head mechanism to split each input token into multiple sub-tokens and distribute them to different experts, MH-MoE achieves denser expert activation without an increase in computational and parameter complexity. Specifically, as shown in Figure 2, when provided with a single input to-

ken, MH-MoE activates four experts by splitting it into four sub-tokens, whereas SMoE only activates one expert. Furthermore, the allocation of sub-tokens to distinct experts enables the model to simultaneously focus on information from various representation spaces within different experts, ensuring a more granular understanding for subtle differences in both vision and language patterns. See in Figure 2, sub-tokens assigned to Experts 3 and 2 capture a detailed understanding of each character's actions within an image patch, while those assigned to Experts 1 and 4 explicitly model the semantics of the false cognate 'camera'. After expert processing, sub-tokens are seamlessly reintegrated into the original token form, thereby circumventing any additional computational burden in subsequent non-parallel layers, e.g., attention layer, while also integrating semantic information captured from multiple experts

MH-MoE maintains following strengths: (1) **Higher experts activation & better scalability**. MH-MoE can alleviate lower expert activation problem and significantly enhance the usage of larger experts by enabling optimization of almost all of experts, e.g., achieving **90.71%** activation in Figure 1 (a), allowing for more efficient scaling of model capacity. (2) **Finer-grained understanding ability**. Multi-head mechanism adopted in MH-MoE assign sub-tokens to different experts, enabling to jointly attend to information from different representation spaces at different experts, and finally achieving better finer-grained understanding ability. For example, refer to the bright area in Figure 1 (b), where sub-tokens are distributed among a more diverse set of experts, facilitating the capture of semantically-rich information. (3) **Seamless integration**. The implementation of MH-MoE is remarkably straightforward and decoupled from other SMoE optimization methods (e.g., GShard (Lepikhin et al., 2020)), making it very easy to integrate them together to achieve better performance.

We evaluate the proposed MH-MoE on three model pre-training and fine-tuning setting: English-focused language modeling, Multi-lingual language modeling and Masked multi-modality modeling. Extensive experimental among

these three tasks demonstrate the effectiveness of MH-MoE.

## 2. Background

**Sparse Mixture of Experts.** Sparse Mixture-of-Experts (SMoE) (Shazeer et al., 2017b; Du et al., 2021; Chi et al., 2022; Clark et al., 2022) enhances model capacity while maintaining a constant computational demand, thus achieving better performance than densely-activated models on various tasks (Lepikhin et al., 2021; Kumatani et al., 2021; Zhao et al., 2023; Pham et al., 2023) and being emerged as a pivotal advancement in the field of deep learning.

Different from densely-activated models, each MoE layer consists of $N$ independent Feed-Forward Networks (FFN) $\{f_i^{\text{FFN}}\}_{i=0}^N$ as the experts, along with a gating function $g(\cdot)$ to model a probability distribution indicating the weights over these experts' outputs. For the hidden representation $\mathbf{h} \in \mathbb{R}^d$ of each input token, the gating value of routing $\mathbf{h}$ to expert $f_i^{\text{FFN}}$ is denoted as:

$$g\left(f_i^{\text{FFN}}\right) = \exp\left(\mathbf{h} \cdot \mathbf{e}_i\right) / \sum_{j=0}^N \exp\left(\mathbf{h} \cdot \mathbf{e}_j\right), \quad (1)$$

where $\mathbf{e}_i$ denotes the trainable embedding of the $i$-th expert and $\sum_{i=0}^N g\left(f_i^{\text{FFN}}\right) = 1$. Then, the corresponding $k$ experts, according to the top-$k$ gated values, are activated and the output $\mathbf{O}$ of the MoE layer is

$$\mathbf{O} = \mathbf{h} + \sum_{i \in \Phi} g\left(f_i^{\text{FFN}}\right) \cdot f_i^{\text{FFN}}\left(\mathbf{h}\right). \quad (2)$$

where $\Phi$ denote activated experts set and $|\Phi| = k$.

**Routing Mechanism in SMoE.** As described above, the most commonly used routing mechanism involves selecting the top-$k$ experts from $N$ experts, where $k \ll N$ (Shazeer et al., 2017a), e.g., $k = 2$ and $N = 2048$ in GShard (Lepikhin et al., 2020). Such a routing mechanism allows the combination of data parallelism and expert parallelism. Yang et al. (2021) and Lepikhin et al. (2020) suggest that larger values of $k$ often contribute to better model performance. However, with the increase in the value of $k$, training models with conventional top-$k$ routing implementation becomes much less efficient (Lepikhin et al., 2020).

In this paper, we introduce MH-MoE, a simple but efficient manner to make denser expert activation without an increase in computational complexity.

## 3. Method

The full architecture of MH-MoE can be seen in Figure 3, MH-MoE addresses low experts activation and confusion over ambiguity of tokens issues by applying a multi-head

mechanism to split each token into sub-tokens and route them to various experts to achieve denser expert activation as well as deeper understanding.

### 3.1. Multi-Head Mixture-of-Experts

Concretely, we denote a sequence of inputs tokens by $\mathbf{X} \in \mathbb{R}^{l \times d}$, where $l$ is the number of tokens and $d$ represents the length of token dimension. In MH-MoE, each parallel layer contains a set of $N$ experts, each presented as $\{f_i^{\text{FFN}} : \mathbb{R}^{\frac{d}{h}} \to \mathbb{R}^{\frac{d}{h}}\}_{i=0}^N$, $h$ denotes the number of heads in multi-head mechanism, which is decoupled from the head in the multi-head self-attention layer. For clarity, we describe the operation of a single MH-MoE layer here only.

First, $\mathbf{X}$ is projected by a multi-head layer with parameter matrices $\mathbf{W}_{\text{head}} \in \mathbb{R}^{d \times d}$,

$$\hat{\mathbf{X}} = \mathbf{X} \cdot \mathbf{W}_{\text{head}}^\top \quad (3)$$

where $\hat{\mathbf{X}} \in \mathbb{R}^{l \times d}$. After that, every token in $\hat{\mathbf{X}}$ is split into $h$ sub-tokens along the token dimensions, and these sub-tokens are arranged in parallel according to the original token sequence, forming a new feature space $\ddot{\mathbf{X}} \in \mathbb{R}^{(l \times h) \times \frac{d}{h}}$ as:

$$
\begin{aligned}
\ddot{\mathbf{X}} &= F_s(\hat{\mathbf{X}}) \\
&= \left[ \overbrace{\underbrace{\mathbf{x}_0^0, \ldots, \mathbf{x}_{h-1}^0, \ldots, \mathbf{x}_j^i, \mathbf{x}_{j+1}^i, \ldots, \mathbf{x}_{h-1}^l}_{l \times h}}^{h} \right],
\end{aligned} \quad (4)
$$

where function $F_s$ denotes the token splitting operation: $\mathbb{R}^{l \times d} \to \mathbb{R}^{(l \times h) \times \frac{d}{h}}$, and each sub-token is presented as $\mathbf{x}_j^i \in \mathbb{R}^{\frac{d}{h}}$, meaning it is the the $j^{th}$ sub-token split from the $i^{th}$ token.

Then all these sub-tokens are fed into the gating function $g(\cdot)$. The gating value of routing a certain sub-token $\mathbf{x}_j^i$ into the $p^{th}$ expert is computed as

$$g\left(f_p^{\text{FFN}}\right) = \frac{\exp\left(\mathbf{x}_j^i \cdot \mathbf{e}_p\right)}{\sum_{\xi=0}^N \exp\left(\mathbf{x}_j^i \cdot \mathbf{e}_\xi\right)}, \quad (5)$$

where $\mathbf{e}_p \in \mathbb{R}^{\frac{d}{h}}$ is the learnable embedding of the $p^{th}$ expert. In this paper, we mainly focus on top-$k$ routing, *i.e.*, only the experts with the largest top-$k$ routing score is activated. $\Phi = \text{Top}_k\left(g\left(f^{\text{FFN}}\right)\right)$ denote the set of activated experts and $|\Phi| = k$. Then $\mathbf{x}_j^i$ is processed by these activated experts as following,

$$\mathbf{o}_j^i = \mathbf{x}_j^i + \sum_{p \in \Phi} g\left(f_p^{\text{FFN}}\right) \cdot f_p^{\text{FFN}}\left(\mathbf{x}_j^i\right). \quad (6)$$

After that, all obtained $\mathbf{o}_j^i$ are rearranged in the original
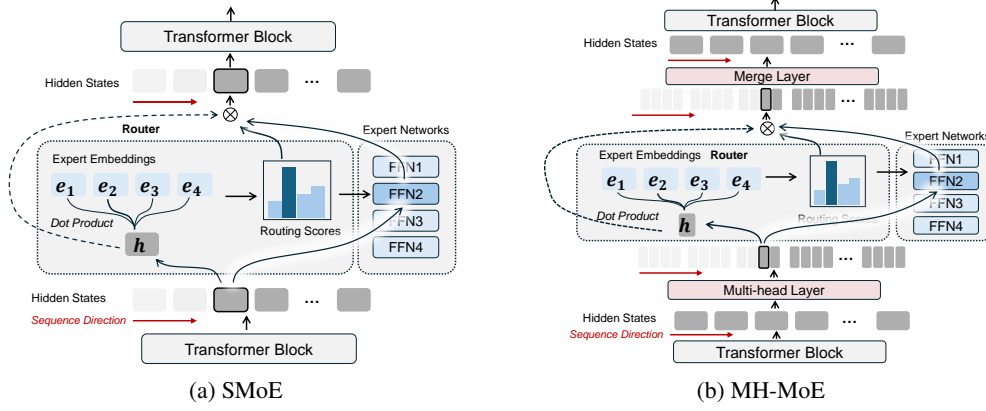
(a) SMoE

(b) MH-MoE

*Figure 3.* **Illustration of a typical SMoE layer and the proposed MH-MoE layer**. (a) An SMoE layer consists of a router and expert networks, where the experts are sparsely activated according to dot-product token-expert routing scores. (b) MH-MoE introduces additional two MLP layers, namely the multi-head layer and merge layer, and a Token-Splitting-Merging (TSM, Eq. 4 and Eq. 8) operation incorporated between these two MLPs.

order of sub-tokens and integrated together as

$$\mathbf{O} = \left[ \overbrace{\mathbf{o}_0^0, \ldots \mathbf{o}_{h-1}^0, \ldots, \mathbf{o}_j^i, \mathbf{o}_{j+1}^i, \ldots, \mathbf{o}_{h-1}^l}^{h} \right], \quad (7)$$

where $\mathbf{O} \in \mathbb{R}^{(l \times h) \times \frac{d}{h}}$. After that, $\mathbf{O}$ is transformed back the into original token form by a token merging operation $\mathcal{F}_{\mathrm{m}}: \mathbb{R}^{(l \times h) \times \frac{d}{h}} \to \mathbb{R}^{l \times d}$:

$$\bar{\mathbf{X}} = \mathcal{F}_{\mathrm{m}} \left( \mathbf{O} \right)^\top, \quad (8)$$

where $\bar{\mathbf{X}} \in \mathbb{R}^{l \times d}$. Finally, $\bar{\mathbf{X}}$ is projected by a merge layer with parameter matrices $\mathbf{W}_{\mathrm{merge}} \in \mathbb{R}^{d \times d}$ to effective integration of multiple features $\mathbf{o}_j^i$ capturing detailed information from different expert representation spaces. The operation is presented as following:

$$\breve{\mathbf{X}} = \bar{\mathbf{X}} \cdot \mathbf{W}_{\mathrm{merge}}^\top. \quad (9)$$

Then we get the final output $\breve{\mathbf{X}}$ of the single MH-MoE layer.

We name the token splitting (Eq. 4) and token merging (Eq. 8) operations together as the Token-Splitting-Mergin (TSM) operation. By implementing the aforementioned operations, we have effectively increased the average volume of data routed to a specific expert by a factor of $h$, as demonstrated in Eq. 4. Consequently, this achievement has resulted in denser expert activation. Furthermore, the allocation of sub-tokens to distinct experts within MH-MoE enables us to collectively capture semantic information from diverse feature spaces across these experts, thereby enhancing the model's ability to achieve a finer-grained understanding.

The operations mentioned above ensure that the shapes of the input and output in the MH-MoE layer remain unchanged. Consequently, no additional computational cost

is introduced in the subsequent block. Specifically, we introduce a hyperparameter $\beta$ to scale the inner dimensions of each expert, aiming to balance the parameters introduced by the multi-head layer and merge layer, aligning the model's parameters and computational complexity with the original SMoE.

As the Pytorch-like style pseudocode of MH-MoE shown in Appendix E, MH-MoE is characterized by its overall simplicity of implementation, necessitating minimal modifications to the SMoE implementation. Additionally, it is decoupled from other SMoE optimization strategies (Lepikhin et al., 2020; Chi et al., 2022), thereby facilitating its convenient integration with other optimized SMoE frameworks to enhance performance.

### 3.2. Training Objectives

The training objective of MH-MoE involves the simultaneous minimization of both the loss associated with the target task and an auxiliary load balancing loss.

**Load balancing loss.** As described in Section 2, there is usually an expert load imbalance problem (Xie et al., 2023; Lepikhin et al., 2020). So, following (Lepikhin et al., 2020; Fedus et al., 2022), given the sub-token set $\ddot{\mathbf{X}}$ (depicted in Eq. 4) and the frequency $t_p$ of how many sub-tokens are routed to the $p^{th}$ expert, we compute the load balancing loss $\mathcal{L}_{\mathrm{balance}}$ via:

$$\mathcal{L}_{\mathrm{balance}} = \frac{N}{|\ddot{\mathbf{X}}|} \sum_{p=1}^{N} \sum_{\mathbf{x}_j^i \in \ddot{\mathbf{X}}} t_p \cdot g \left( f_p^{\mathrm{FFN}} \right), \quad (10)$$

where $N$ denotes the number of experts, $|\ddot{\mathbf{X}}|$ is the number of sub-tokens contained in $\ddot{\mathbf{X}}$. $g \left( f_p^{\mathrm{FFN}} \right)$ is the gating function depicted in Eq. 5, denoting the gating value of routing a certain sub-token $\mathbf{x}_j^i$ into the $p^{th}$ expert.
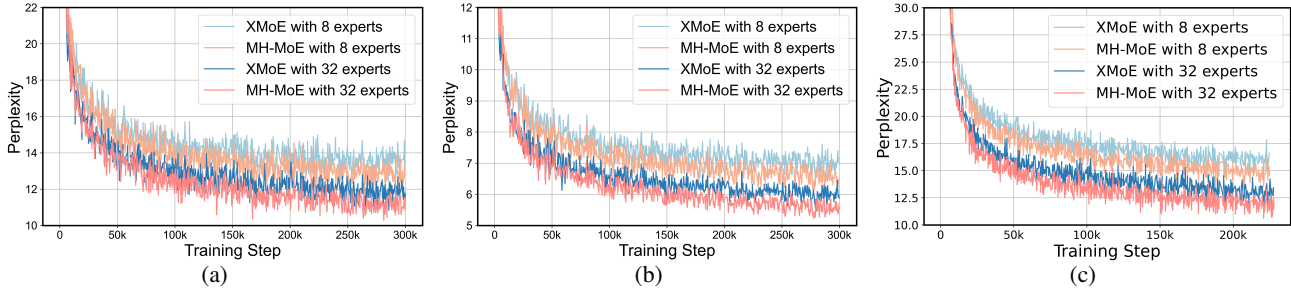
*Figure 4.* **Perplexity on validation dataset during the training phase** reported for Dense, X-MoE and MH-MoE across three pre-training tasks. (a) *English-focused language modeling.* (b) *Multi-lingual language modeling.* (c) *Masked multi-modal modeling*

**Task specific loss.** The term $\mathcal{L}_{task}$ is dependent on the particular task that MH-MoE is designed to learn. For instance, during pre-training, we utilize the language modeling loss (Radford et al., 2018), whereas the model predicts the next word in a sequence.

So, the overall training objective is to minimize:

$$\mathcal{L} = \mathcal{L}_{task} + \alpha\mathcal{L}_{balance}, \qquad (11)$$

where $\alpha$ is a coefficient for load balancing.

## 4. Experiments

### 4.1. Experimental Setup

**Compared Baselines.** We include two baseline models for comparison purposes: (1) **Dense**, which represents a Transformer decoder without the incorporation of sparsely-activated parallel modules (i.e., SMoE layer). (2) **X-MoE**, which is our implementation based on the approach proposed by Chi et al. (2022). We build our MH-MoE upon X-MoE and uses identical settings to those employed in X-MoE. Note that the all these models are pre-trained using the same training data as MH-MoE, and we ensure that the parameter count of our model remains consistent with or lower than that of X-MoE, ensuring a fair and equitable comparison. A detailed analysis and comparison about parameter and computational complexity can be found in Section 5.3 and Table 11.

**Pre-training Data.** We detail the pre-training data of MH-MoE, demonstrating its effectiveness in enabling denser expert activation and finer-grained understanding through a series of experiments. These experiments are organized into three thematic categories: (1) For the English-focused experiments, we pretrain both the baseline models and MH-MoE on the RedPajama dataset (Computer, 2023), which is an open-source pre-training dataset comprising sources such as Common Crawl, C4 (Raffel et al., 2020), Wikipedia, and additional curated datasets. The pretraining is conducted using GPT tasks to predict the next word in a sequence. (2) In the context of multilingual representation, we pretrain the baseline models and MH-MoE on the multilingual Wikipedia, following the approach described

in XLM (Lample & Conneau, 2019), again utilizing GPT tasks. (3) For the multimodal domain, we pretrain all compared baselines and MH-MoE on masked multi-modality modeling task on both monomodal and multimodal data (14M images, 160GB documents and 21M image-text pairs following Wang et al. (2022)), and we present the details of these pre-training data in Appendix A.

**Model Architecture and Hyperparameters.** For all experiments, we use the X-MoE Chi et al. (2022) as our backbone architecture to build our MH-MoE, which has shown better performance than prior SMoE models such as Switch Transformers (Fedus et al., 2022) on cross-lingual understanding benchmarks. For English-focused Language Modeling and Multi-lingual Language Modeling, we construct Dense, X-MoE and MH-MoE using the Transformer (Vaswani et al., 2017) decoder (L = 12, H = 768, A = 12) with the GPT-4[1] vocabulary as the backbone architecture. The pre-training procedure takes 14 days on 2 NVIDIA DGX-2 Stations. For Masked Multi-modal Modeling, we build Dense, X-MoE and MH-MoE following the same Transformer encoder architecture as BEiT v3 (Wang et al., 2022). The pre-training procedure takes 4 days on 2 NVIDIA DGX-2 Stations. For all three pre-training tasks, we set the head number $h = 4$. More details about architecture and training hyperparameters can be found in Appendix B and C.

### 4.2. Perplexity Evaluation

We examined the validation perplexity curves for all pre-trained models and pre-training tasks under two expert settings (8 experts and 32 experts). The perplexity trends are depicted in Figure 4, with the final perplexity values listed in Table 1. We can observe that as training progresses: 1) the perplexity of our MH-MoE remained lower in comparison to the compared baselines, indicating more effective learning; 2) MH-MoE achieved the lowest perplexity across three distinct experimental setups; 3) an increase in the number of experts led to a corresponding decrease in the perplexity of MH-MoE, suggesting that the model ben-

---

[1] https://github.com/openai/tiktoken

Table 1. Results of upstream perplexity evaluation. We report the validation perplexity cross two setting: 8 experts and 32 experts.

| Model | Perplexity ↓ | |
|---|---|---|
| | 8 Experts | 32 Experts |
| *English-focused language modeling* | | |
| Dense (without Experts) | 16.23 | 16.23 |
| X-MoE | 14.82 | 11.96 |
| MH-MoE (Ours) | **12.72** | **10.28** |
| *Multi-lingual language modeling* | | |
| Dense (without Experts) | 8.56 | 8.56 |
| X-MoE | 7.19 | 6.02 |
| MH-MoE (Ours) | **6.26** | **5.09** |
| *Masked multi-modal modeling* | | |
| Dense (without Experts) | 17.95 | 17.95 |
| X-MoE | 16.34 | 12.68 |
| MH-MoE (Ours) | **14.73** | **10.87** |

efits from enhanced representation learning capabilities as more experts are incorporated. These results collectively demonstrate the superiority of MH-MoE in terms of learning efficiency and language representation across multiple pre-training paradigms.

### 4.3. Downstream Evaluation

For each pre-training task, we conduct corresponding downstream evaluation to validate the efficacy of MH-MoE.

**English-focused Language Modeling.** We evaluate our models on a total of 9 different zero-shot benchmarks to assess their abilities across various natural language tasks like common sense reasoning, general language understanding and knowledge understanding using the LLM Evaluation Harness (Gao et al., 2023). As shown in Table 2, comparing X-MoE with the Dense model, X-MoE show notable improvement, indicating that SMoE models (e.g., X-MoE) benefit from the large model capacity. Overall, for all benchmarks, our MH-MoE obtains the best performance, achieving an average performance gain of 1.1 for 8 experts setting and 1.5 for 32 experts setting compared to X-MoE, demonstrating the effectiveness of our proposed multi-head mechanism on modeling English-focused language.

**Multi-lingual Language Modeling.** We evaluate our multi-lingual language models on the cross-lingual natural language inference (XNLI) corpus (Conneau et al., 2018), which is the extension of the multi-genre NLI (MultiNLI) corpus to 14 languages. We follow the LLM Evaluation Harness pipeline and use the zero-shot setting to evaluate the multi-lingual ability. Table 3 presents the zero-shot evaluation results on XNLI task. Similarly, X-MoE benefit from the large model capacity and show notable improvement compared with Dense model. Overall, MH-MoE obtains the best performance, surpassing X-MoE by an average performance gain of 0.6 for 8 experts setting and 0.8 for

32 experts setting. Comparing MH-MoE with the X-MoE, it shows that MH-MoE models provide consistent gains on downstream tasks, demonstrating the effectiveness of our proposed multi-head mechanism on modeling cross-lingual natural language.

**Masked Multi-modal Modeling.** We evaluate on the widely used vision-language understanding and generation benchmarks, including visual question answering (Goyal et al., 2017), visual reasoning (Suhr et al., 2019) and image captioning (Lin et al., 2014). We report *vqa-score* on VQAv2, accuracy for NLVR2. For COCO image captioning, we report BLEU@4 (B@4), METEOR (M), CIDEr (C), and SPICE (S). Table 4 presents the evaluation results. For VQA task, MH-MoE outperforms both Dense and X-MoE by a large margin, e.g., 4.24 and 1.69 points gain on test-dev split, respectively. For visual reasoning task, MH-MoE beats all these two baselines on both dev (1.5 points gain than X-MoE) and test-P split (1.7 points gain than X-MoE). For image captioning task, MH-MoE surpasses X-MoE by 4.2%, 10.2%, 9.4% in terms of B@4, M and S, respectively. Above results show that X-MoE exhibits enhanced visual information comprehension, which also validates the effectiveness of our proposed multi-head mechanism in capturing diverse semantic and detailed information within visual data.

### 4.4. Ablation Studies

This section presents experimental analysis to demonstrate the functionality of MH-MoE. In all comparative experiments, *we ensure parameter equality across models by adjusting the internal dimensions of the experts*.

**Number of Heads $h$.** We conduct experiments by adjusting the number of heads ($h$ = 2, 4, 6, 8, and 12) in MH-MoE. As shown in Table 5, we find that across all settings of $h$, our model consistently outperforms the X-MoE, demonstrating the effectiveness of MH-MoE. Besides, as the value of $h$ increases, we observe an initial improvement followed by a decline in our model's performance. This leads us to hypothesize that when $h \leq 6$ the enhancement in performance benefits from the multi-head mechanism by activating a greater number of experts, thereby enhancing the model's effectiveness and capturing a wider range of fine-grained token information. However, as $h$ continues to increase beyond 6, the excessive subdivision of tokens may inadvertently impair their original semantic content, resulting in a decrease in model performance.

**Effect of MH-MoE Components.** As shown in Figure 3 (b), the multi-head mechanism utilized in our MH-MoE primarily incorporates two components: the Multilayer Perceptron (MLP) layers, including the multi-head layer (Eq. 3) and merge layer (Eq. 9), and the Token-Splitting-Merging (TSM) operation (Eq. 4 and Eq. 8). We conduct

*Table 2.* Accuracy / accuracy-normalization scores for language understanding tasks using the LLM Evaluation Harness (Gao et al., 2023).
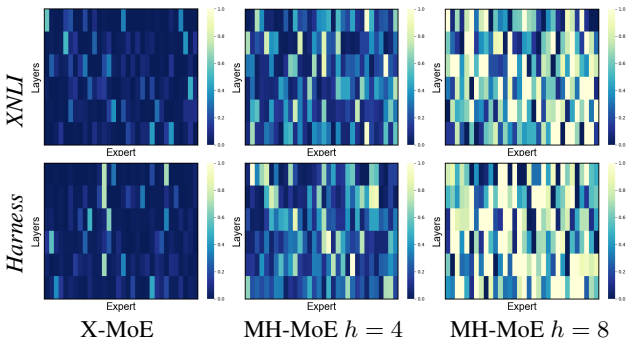
| Model | ARC-Challenge | ARC-Easy | RTE | BookQA | Winogrande | PiQA | BoolQ | HellaSwag | TruthfulQA (mc1/mc2) | Avg |
|---|---|---|---|---|---|---|---|---|---|---|
| Dense | 18.1/23.3 | 44.9/39.7 | 51.5 | 17.1/29.0 | 48.2 | 66.6 | 55.0 | 29.7/34.1 | 24.1/39.3 | 37.2 |
| *Experts Number $N = 8$* | | | | | | | | | | |
| X-MoE | 19.0/24.7 | 48.3/42.0 | 52.7 | 17.4/29.8 | 50.3 | 67.9 | 58.4 | 31.4/35.7 | 24.3/40.2 | 38.7 |
| MH-MoE | **19.6/25.2** | **50.2/42.2** | **53.0** | **18.2/30.3** | **51.1** | **68.7** | **59.6** | **33.2/40.3** | **24.7/40.9** | **39.8** |
| *Experts Number $N = 32$* | | | | | | | | | | |
| X-MoE | 19.4/24.8 | 50.4/42.5 | 52.7 | 17.8/30.0 | 51.3 | 68.8 | 52.8 | 33.4/40.1 | 24.3/39.1 | 39.1 |
| MH-MoE | **21.4/26.8** | **50.6/44.8** | **53.4** | **18.8/31.6** | **53.8** | **69.3** | **56.6** | **35.0/42.1** | **24.8/39.5** | **40.6** |

*Table 3.* Accuracy / accuracy-normalization scores on multilingual understanding tasks using the LLM Evaluation Harness (Gao et al., 2023).

| Model | bg | de | el | en | es | fr | hi | ru | sw | th | tr | ur | vi | zh | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Dense | 33.1 | 33.3 | 33.0 | 35.1 | 32.8 | 34.4 | 33.6 | 34.2 | 33.3 | 33.1 | 33.3 | 33.9 | 33.5 | 32.9 | 33.5 |
| *Experts Number $N = 8$* | | | | | | | | | | | | | | | |
| X-MoE | 33.9 | **33.4** | 33.4 | 37.3 | 33.3 | 35.9 | 34.5 | 35.0 | 33.5 | 33.6 | 33.4 | 34.2 | 33.3 | 33.2 | 34.1 |
| MH-MoE | **34.4** | 33.2 | **33.9** | **40.1** | **34.0** | **36.4** | **34.6** | **35.2** | **33.8** | **34.4** | 33.3 | **34.7** | **34.6** | **33.5** | **34.7** |
| *Experts Number $N = 32$* | | | | | | | | | | | | | | | |
| X-MoE | 34.5 | 34.5 | 33.4 | 39.6 | 33.1 | 35.3 | 34.1 | 35.4 | 33.6 | 34.7 | 33.7 | 33.6 | 34.5 | 33.3 | 34.5 |
| MH-MoE | **35.8** | **35.6** | **34.1** | **40.7** | **33.9** | **36.7** | **34.4** | **36.3** | **34.3** | **36.0** | **34.1** | **34.3** | **35.2** | **33.6** | **35.3** |

*Table 4.* Results of visual question answering, visual reasoning, and image captioning tasks.

| Model | VQAv2 | | NLVR2 | | COCO Captioning | | | |
|---|---|---|---|---|---|---|---|---|
| | test-dev | test-std | dev | test-P | B@4 | M | C | S |
| Dense | 65.9 | 66.0 | 73.8 | 74.2 | 35.9 | 29.3 | 120.5 | 19.6 |
| *Experts Number $N = 8$* | | | | | | | | |
| X-MoE | 68.4 | 69.7 | 75.5 | 76.1 | 38.1 | 30.2 | 122.9 | 21.3 |
| MH-MoE | **70.1** | **71.4** | **77.0** | **77.8** | **39.7** | **33.1** | **124.1** | **23.0** |



*Figure 5.* **Distribution of expert activation in X-MoE and MH-MoE** on both *Harness* (Gao et al., 2023) and *XNLI* (Conneau et al., 2018) corpus, encompassing 6 SMoE layers with 32 experts per layer. The top of the heatmap is the first SMoE layer while the bottom is the last. Experts activation ratio is determined by calculating the ratio of each expert's selection frequency in each MoE layer to the total number of tokens.

a detailed analysis of the effectiveness of each component within our model, as well as the necessity of their integration.

The results are presented in Table 6. A comparative analysis between Dense v.s. Dense$_{w/o \text{ MLP}}$, as well as X-MoE v.s. X-MoE$_{w/ \text{ MLP}}$, reveals that introduction of the MLP layer does not enhance the model's performance. Similarly, when comparing MH-MoE with MH-MoE$_{w/o \text{ MLP}}$, it becomes evident that the inclusion of only the MLP, in the absence of the TS, also does not yield any improvement in the model's effectiveness. The parameter quantities of the models being compared pairwise are equal.

An intriguing observation is made when comparing MH-MoE with MH-MoE$_{w/o \text{ TS}}$. Introducing Token-Splitting-Merging (TSM) alone, without MLP, results in a slight in-

crease in model performance. In contrast, a significant enhancement in model performance is only achieved when both MLP and TS are incorporated simultaneously. We hypothesize that introduction of TS, without the integration of MLP, activates more experts, but the segmentation and merging of the model appears overly straightforward and abrupt in its execution. This limitation hinders the model's ability to meaningfully segment tokens into sub-tokens and effectively merge the diverse information gathered from different expert spaces.

**Number of MLP layers.** We explore the impact of varying the number of layers ($n = 0, 1, 2, 3$) in MLP on MH-MoE performance. For configurations exceeding a single layer, ReLU activation functions were incorporated between MLP layers to ensure the non-linearity of transformations. The parameter quantities of the models being compared are equal. Upon analyzing the results in Table 7, we observe that increasing the number of MLP layers beyond one had a negligible impact on the model's performance. This indicates that a single-layer MLP is sufficient for accomplishing token segmentation and fusion.

## 5. Analysis

### 5.1. Experts Activation Analysis

**Experts Activation.** We visualize the activation of each expert varies across parallel expert layers for X-MoE and MH-MoE at Figure 5. It can be observed that: 1) X-MoE demonstrate a more skewed distribution, wherein a significant portion of experts remain inactivated all the time. 2) Our MH-MoE achieves a denser expert activation compared to X-MoE, effectively mitigating the issue of low expert utilization. 3) As the number of heads $h$ increases, the

*Table 5.* Comparison results for different head number $h$. S-Dim denotes the dimension length of sub-tokens.

| Model | Heads $h$ | S-Dim | Perplexity |
|---|---|---|---|
| X-MoE | - | - | 14.82 |
| | 2 | 384 | 12.87 |
| | 4 | 192 | 12.72 |
| MH-MoE | 6 | 128 | **12.41** |
| | 8 | 96 | 12.95 |
| | 12 | 64 | 13.28 |

*Table 6.* Ablation studies of MH-MoE components: MLP layers and the Token-Splitting-Merging (TSM, Eq. 4 and Eq. 8) operation.

| Model | MLP | TSM | Perplexity |
|---|---|---|---|
| Dense | ✗ | ✗ | 16.23 |
| Dense$_{w/ \text{ MLP}}$ | ✓ | ✗ | 16.40 |
| X-MoE | ✗ | ✗ | 14.82 |
| X-MoE$_{w/ \text{ MLP}}$ | ✓ | ✗ | 14.77 |
| MH-MoE$_{w/o \text{ TS}}$ | ✓ | ✗ | 14.77 |
| MH-MoE$_{w/o \text{ MLP}}$ | ✗ | ✓ | 13.97 |
| MH-MoE | ✓ | ✓ | **12.72** |

*Table 7.* Comparison results for different numbers of MLP layers $n$. The results are averaged over five runs.

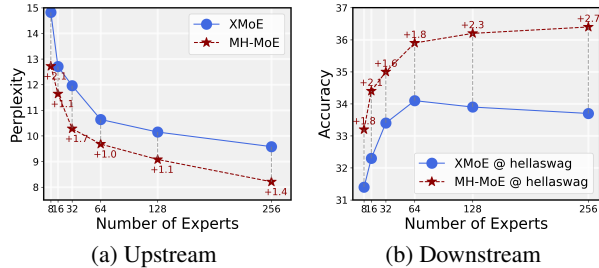| $n$ | Upstream | Downstream | | |
|---|---|---|---|---|
| | Perplexity | RTE | PIQA | Winogrande |
| 0 | 13.97 | 52.9 | 68.2 | 51.7 |
| 1 | 12.72 | 53.4 | **69.3** | **53.8** |
| 2 | **12.66** | **54.0** | 68.8 | 53.3 |
| 3 | 12.87 | 53.1 | 68.8 | 52.7 |



(a) Upstream  (b) Downstream

*Figure 6.* **Upstream and downstream results for scaling up the number of experts in X-MoE and MH-MoE**. (a) Training perplexity ($\downarrow$) when scaling the number of experts. (b) Downstream performance (accuracy scores $\uparrow$) on hellaswag when scaling the number of experts.
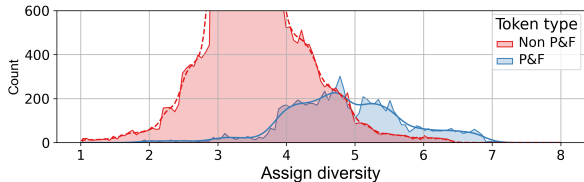


*Figure 7.* **Comparison for sub-tokens assign diversity** (the number of different experts they are routed to) for P&F and Non P&F tokens. P&F tokens refer to the polysemous and false cognate words identified by GPT-4, while Non P&F tokens represent the remaining words.

expert activation frequency in MH-MoE also rises.

**Scalability.** We explore the scalability for both X-MoE and MH-MoE by scaling up the number of experts from 8 to 256 (about 7B parameters). For upstream performance, as shown in Figure 6 (a), with the increase of experts, our MH-MoE could bring more gains. It is because MH-MoE could mitigate the low expert activation problem effectively. With this ability, the superiority of the large-scale SMoE model will be better exerted, thereby achieving the improvement of the upper bound of SMoE with more experts. Detailed validation perplexity curves for these scaling up experiments can be found in Figure 9 at Appendix F. For downstream performance shown in Figure 6 (b), for X-MoE, expert number = 64 is the upper bound, meaning that continuing to increase the number of experts will not bring any gain. Our MH-MoE not only has a performance advantage over the X-MoE with the same number of experts, but also improves the upper bound from 64 to 256, which demonstrates the effectiveness of the scalability of our MH-MoE
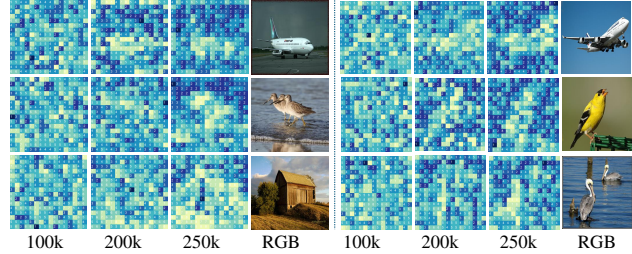


*Figure 8.* **Assign diversity of sub-tokens split from different patches** in vision data with respect to training steps (100k $\rightarrow$ 200k $\rightarrow$ 250k steps). Brighter regions indicate sub-tokens split from this patch are distributed to a greater number of diverse experts.

on downstream tasks.

### 5.2. Fine-grained understanding Analysis

In Section 4, our model excels in multiple upstream and downstream tasks, demonstrating superior fine-grained modeling capabilities, both for languages and images. In this section, we delve into a more granular analysis to validate how the multi-head mechanism aids MH-MoE in capturing diverse and intricate semantic information that is often challenging to comprehend, e.g., polysemous and false cognates words (denoted as PF tokens) in languages, and semantically-rich areas in images. Note that for languages data, we utilized the GPT-4 API (OpenAI, 2023) to extract polysemous words and false cognates from the XNLI (Conneau et al., 2018) corpus, and the corresponding prompt can be found in Table 12.

**Experts Assign within Token.** For languages data, we compute and compare the divergence levels (i.e., the number of different experts these sub-tokens are routed to) of sub-tokens split from PF tokens and Non-PF tokens. We conduct on MH-MoE with 8 heads ($h$=8) and represent the divergence of each token by calculating the mean divergence across the model's various layers. The results, presented in Figure 7, clearly demonstrate that the distribution of divergence for PF tokens is significantly skewed towards the right when compared to that of Non-PF tokens. This indicates that, in the MH-MoE's inference process, PF tokens route their sub-tokens to a greater number of different experts, thereby capturing diverse semantic information in

contrast to Non-PF tokens for a better polysemous and false cognates word modeling.

For image data, we analyzed how the divergence levels of different patches evolve during the training process, as illustrated in Figure 8. Interestingly, we observe that as the training steps increase, the divergence levels gradually increase in high-frequency texture regions (or regions with rich semantics), while the divergence levels in low-frequency texture regions gradually decrease. This indicates that during the training process, MH-MoE tends to route tokens from areas with complex textures to a greater variety of experts, thereby enhancing the finer-grained understanding of the semantics in that region. For more visualization examples, please refer to the Figure 10 at Appendix G.

### 5.3. Complexity & Parameter Analysis.

We present a analysis of Complexity & Parameter for X-MoE and MH-MoE in Appendix D, to validate that for all experiments setting, the computational and parameter cost of our MH-MoE are both lower than SMoE. Besides, a detailed parameter count for all experiments and comparable models can be seen in Table 11.

## 6. Conclusion

In this paper, we study how we can to achieve a denser experts activation without introducing additional cost, while improving the fine-grained understanding ability. With the proposed Multi-Head Mixture-of-Experts, we can easily implement the aforementioned functionality. Furthermore, the simplicity of MH-MoE allows it to integrate with other SMoE frameworks to enhance performance easily. Extensive empirical results across three tasks demonstrate the effectiveness of MH-MoE.

## 7. Broader Impact

In previous NLP pipelines, the dimension of word tokens has been conventionally maintained constant during both training and inference stages. We are the first to attempt token segmentation outside of the multi-head attention module, aiming to enhance the model's capabilities in several respects, including a more nuanced and multifaceted understanding of the token content as well as fostering a sparser network architecture. We believe this to be a counterintuitive yet worthwhile exploration in the field.

## References

Bao, H., Dong, L., Wei, F., Wang, W., Yang, N., Liu, X., Wang, Y., Gao, J., Piao, S., Zhou, M., and Hon, H. UniLMv2: Pseudo-masked language models for unified language model pre-training. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 642–652. PMLR, 2020. URL http://proceedings.mlr.press/v119/bao20a.html.

Changpinyo, S., Sharma, P., Ding, N., and Soricut, R. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pp. 3558–3568. Computer Vision Foundation / IEEE, 2021.

Chi, Z., Dong, L., Huang, S., Dai, D., Ma, S., Patra, B., Singhal, S., Bajaj, P., Song, X., Mao, X.-L., et al. On the representation collapse of sparse mixture of experts. *Advances in Neural Information Processing Systems*, 35: 34600–34613, 2022.

Chung, H. W., Hou, L., Longpre, S., Zoph, B., Tay, Y., Fedus, W., Li, Y., Wang, X., Dehghani, M., Brahma, S., et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.

Clark, A., Casas, D. d. l., Guy, A., Mensch, A., Paganini, M., Hoffmann, J., Damoc, B., Hechtman, B., Cai, T., Borgeaud, S., et al. Unified scaling laws for routed language models. *arXiv preprint arXiv:2202.01169*, 2022.

Computer, T. Redpajama: An open source recipe to reproduce llama training dataset, 2023. URL https://github.com/togethercomputer/RedPajama-Data.

Conneau, A., Rinott, R., Lample, G., Williams, A., Bowman, S., Schwenk, H., and Stoyanov, V. XNLI: Evaluating cross-lingual sentence representations. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2475–2485, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1269. URL https://www.aclweb.org/anthology/D18-1269.

Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., and Stoyanov, V. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020*, pp. 8440–8451. Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.acl-main.747. URL https://doi.org/10.18653/v1/2020.acl-main.747.

Du, N., Huang, Y., Dai, A. M., Tong, S., Lepikhin, D., Xu, Y., Krikun, M., Zhou, Y., Yu, A. W., Firat, O., et al. Glam: Efficient scaling of language models with

mixture-of-experts. *arXiv preprint arXiv:2112.06905*, 2021.

Fedus, W., Zoph, B., and Shazeer, N. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *The Journal of Machine Learning Research*, 23(1):5232–5270, 2022.

Gao, L., Tow, J., Abbasi, B., Biderman, S., Black, S., DiPofi, A., Foster, C., Golding, L., Hsu, J., Le Noac'h, A., Li, H., McDonell, K., Muennighoff, N., Ociepa, C., Phang, J., Reynolds, L., Schoelkopf, H., Skowron, A., Sutawika, L., Tang, E., Thite, A., Wang, B., Wang, K., and Zou, A. A framework for few-shot language model evaluation, 12 2023. URL https://zenodo.org/records/10256836.

Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., and Parikh, D. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pp. 6325–6334. IEEE Computer Society, 2017. doi: 10.1109/CVPR.2017.670. URL https://doi.org/10.1109/CVPR.2017.670.

Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D. S., Casas, D. d. l., Hanna, E. B., Bressand, F., et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.

Krishna, R., Zhu, Y., Groth, O., Johnson, J., Hata, K., Kravitz, J., Chen, S., Kalantidis, Y., Li, L., Shamma, D. A., Bernstein, M. S., and Fei-Fei, L. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *Int. J. Comput. Vis.*, 123(1): 32–73, 2017.

Kudo, T. and Richardson, J. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 66–71, Brussels, Belgium, November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-2012. URL https://www.aclweb.org/anthology/D18-2012.

Kumatani, K., Gmyr, R., Salinas, F. C., Liu, L., Zuo, W., Patel, D., Sun, E., and Shi, Y. Building a great multi-lingual teacher with sparsely-gated mixture of experts for speech recognition. *arXiv preprint arXiv:2112.05820*, 2021.

Lample, G. and Conneau, A. Cross-lingual language model pretraining, 2019.

Lepikhin, D., Lee, H., Xu, Y., Chen, D., Firat, O., Huang, Y., Krikun, M., Shazeer, N., and Chen, Z. Gshard: Scaling giant models with conditional computation and automatic sharding. *arXiv preprint arXiv:2006.16668*, 2020.

Lepikhin, D., Lee, H., Xu, Y., Chen, D., Firat, O., Huang, Y., Krikun, M., Shazeer, N., and Chen, Z. {GS}hard: Scaling giant models with conditional computation and automatic sharding. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=qrwe7XHTmYb.

Lin, T., Maire, M., Belongie, S. J., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft COCO: common objects in context. In Fleet, D. J., Pajdla, T., Schiele, B., and Tuytelaars, T. (eds.), *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V*, volume 8693 of *Lecture Notes in Computer Science*, pp. 740–755. Springer, 2014.

Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692, 2019. URL http://arxiv.org/abs/1907.11692.

OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023. URL https://api.semanticscholar.org/CorpusID:257532815.

Ordonez, V., Kulkarni, G., and Berg, T. L. Im2text: Describing images using 1 million captioned photographs. In Shawe-Taylor, J., Zemel, R. S., Bartlett, P. L., Pereira, F. C. N., and Weinberger, K. Q. (eds.), *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain*, pp. 1143–1151, 2011. URL https://proceedings.neurips.cc/paper/2011/hash/5dd9db5e033da9c6fb5ba83c7a7ebea9-Abstract.html.

Peng, Z., Wang, W., Dong, L., Hao, Y., Huang, S., Ma, S., and Wei, F. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.

Pham, H., Kim, Y. J., Mukherjee, S., Woodruff, D. P., Poczos, B., and Awadalla, H. H. Task-based moe for multitask multilingual machine translation. *arXiv preprint arXiv:2308.15772*, 2023.

Radford, A., Narasimhan, K., Salimans, T., and Sutskever, I. Improving language understanding by generative pretraining. 2018. URL http://openai-assets.

s3.amazonaws.com/research-covers/ language-unsupervised/language_ understanding_paper.pdf.

Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL http://jmlr. org/papers/v21/20-074.html.

Sharma, P., Ding, N., Goodman, S., and Soricut, R. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In Gurevych, I. and Miyao, Y. (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, pp. 2556–2565. Association for Computational Linguistics, 2018. URL https: //aclanthology.org/P18-1238/.

Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., and Dean, J. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017a.

Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q., Hinton, G., and Dean, J. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations*, 2017b. URL https://openreview.net/ forum?id=B1ckMDqlg.

Suhr, A., Zhou, S., Zhang, A., Zhang, I., Bai, H., and Artzi, Y. A corpus for reasoning about natural language grounded in photographs. In Korhonen, A., Traum, D. R., and Màrquez, L. (eds.), *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pp. 6418–6428. Association for Computational Linguistics, 2019. doi: 10.18653/v1/p19-1644. URL https://doi.org/ 10.18653/v1/p19-1644.

Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

Trinh, T. H. and Le, Q. V. A simple method for commonsense reasoning. *ArXiv*, abs/1806.02847, 2018.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, pp. 5998–6008, 2017. URL https://proceedings. neurips.cc/paper/2017/hash/ 3f5ee243547dee91fbd053c1c4a845aa-Abstract. html.

Wang, W., Bao, H., Dong, L., Bjorck, J., Peng, Z., Liu, Q., Aggarwal, K., Mohammed, O. K., Singhal, S., Som, S., et al. Image as a foreign language: Beit pretraining for all vision and vision-language tasks. *arXiv preprint arXiv:2208.10442*, 2022.

Xie, Y., Huang, S., Chen, T., and Wei, F. Moec: Mixture of expert clusters. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 13807–13815, 2023.

Yang, A., Lin, J., Men, R., Zhou, C., Jiang, L., Jia, X., Wang, A., Zhang, J., Wang, J., Li, Y., et al. M6-t: Exploring sparse expert models and beyond. *arXiv preprint arXiv:2105.15082*, 2021.

Zhao, X., Chen, X., Cheng, Y., and Chen, T. Sparse moe with language guided routing for multilingual machine translation. In *Conference on Parsimony and Learning (Recent Spotlight Track)*, 2023.

Zhu, Y., Kiros, R., Zemel, R., Salakhutdinov, R., Urtasun, R., Torralba, A., and Fidler, S. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pp. 19–27, 2015.

## A. Pre-training Data of Masked multi-modal modeling task

Table 8 presents of pre-training data in Masked multi-modal modeling task. For multi-modal data, there are about 15M images and 21M image-text pairs collected from five public datasets: Conceptual 12M (CC12M) (Changpinyo et al., 2021), Conceptual Captions (CC3M) (Sharma et al., 2018), SBU Captions (SBU) (Ordonez et al., 2011), COCO (Lin et al., 2014) and Visual Genome (VG) (Krishna et al., 2017). For monomodal data, we use 14M images from ImageNet-21K and 160GB text corpora (Bao et al., 2020) from English Wikipedia, BookCorpus (Zhu et al., 2015), OpenWebText[2], CC-News (Liu et al., 2019), and Stories (Trinh & Le, 2018).

S

*Table 8.* Pretraining data of Masked multi-modal modeling task. All the data are academically accessible.

| Data | Source | Size |
|---|---|---|
| Image-Text Pair | CC12M, CC3M, SBU, COCO, VG | 21M pairs |
| Image | ImageNet-21K | 14M images |
| Text | English Wikipedia, BookCorpus, OpenWebText, CC-News, Stories | 160GB documents |

## B. Model Hyperparameters of Language modeling tasks

Table 9 presents the model hyperparameters of X-MoE and MH-MoE for both English-focused language modeling and Multi-lingual language modeling tasks. The gating temperature $\tau_0$ is initialized as 0.3 and 0.07 for the softmax gating and sigmoid gating, respectively. We use the same vocabulary as XLM-R (Conneau et al., 2020) with 250K subwords tokenized by SentencePiece (Kudo & Richardson, 2018).

*Table 9.* Model hyperparameters of Dense, X-MoE and MH-MoE. The SMoE frequency refers to how many experts each token will be assigned to, i.e., the value of k in the Top- expert selection.

| Hyperparameters | Dense | X-MoE | MH-MoE |
|---|---|---|---|
| FFNs within layer | 2 | 2 | 2 |
| Expert embedding dimension | - | 16 | $16/h$ |
| Initialized gating temperature $\tau_0$ | - | 0.3 / 0.07 | 0.3 / 0.07 |
| Transformer blocks | 12 | 12 | 12 |
| Hidden size | 768 | 768 | 768 |
| FFN inner hidden size | 3,072 | 3,072 | $3,072 \times \beta$ |
| Attention heads | 12 | 12 | 12 |
| SMoE frequency | - | 2 | 2 |

## C. Hyperparameters for Pre-training

Table 10 presents the hyperparameters for pre-training across three tasks: Language modeling tasks (English-focused language modeling and Multi-lingual language modeling tasks) and Masked multi-modal modeling task.

---

[2]http://skylion007.github.io/OpenWebTextCorpus

*Table 10.* Pre-training hyperparameters for Language modeling tasks (English-focused language modeling and Multi-lingual language modeling tasks) and Masked multi-modal modeling task tasks.

| Hyperparameters | Language modeling tasks | Multi-modality modeling task |
|---|---|---|
| Batch size | 256 | 512 |
| Optimizer | Adam | AdamW |
| Batch tokens per task | 1M | - |
| Adam $\epsilon$ | 1e-6 | 1e-6 |
| Adam $\beta$ | (0.9, 0.98) | (0.9, 0.98) |
| Maximum learning rate | 5e-4 | 2.8e-3 |
| Learning rate schedule | Linear decay | Cosine decay |
| Warmup steps | 10,000 | 10,000 |
| Weight decay | 0.01 | 0.05 |
| Transformer dropout | 0.1 | 0.1 |
| Dropout | 0 | 0 |
| Load balancing coefficient | 1e-2 | 1e-2 |

*Table 11.* Parameter count setting of X-MoE and MH-MoE in our experiments for English-focused language modeling, Multi-lingual language modeling and Masked multi-modality modeling tasks. "non-expert param" refers to the parameters that are not part of the expert networks, such as the attention layer, router, etc., while "expert params" represents the total number of parameters in the parallel expert networks. For Dense models, since there are no expert network layers, we only list the total number of parameters. All models under the same task utilize the same architecture and hyperparameters, following identical training settings and steps.

| Expert Setting | Dense | X-MoE | | | MH-MoE | | |
|---|---|---|---|---|---|---|---|
| | **Sum** | non-expert params | expert params | **Sum** | non-expert params | expert params | **Sum** |
| *English-focused language modeling* | | | | | | | |
| 0 expert | 162M | - | - | - | - | - | - |
| 8 experts | - | 134M | 227M | 361M | 141M | 213M | 354M |
| 16 experts | - | 134M | 454M | 588M | 141M | 430M | 571M |
| 32 experts | - | 134M | 908M | 1,042M | 141M | 898M | 1,039M |
| 64 experts | - | 134M | 1,815M | 1,949M | 141M | 1,797M | 1,938M |
| 128 experts | - | 134M | 3,631M | 3,765M | 141M | 3,624M | 3,765M |
| 256 experts | - | 134M | 7,263M | 7,397M | 141M | 7,230M | 7,371M |
| *Multi-lingual language modeling* | | | | | | | |
| 0 expert | 162M | - | - | - | - | - | - |
| 8 experts | - | 134M | 227M | 361M | 141M | 213M | 354M |
| 32 experts | - | 134M | 908M | 1,042M | 141M | 898M | 1,039M |
| *Masked multi-modality modeling* | | | | | | | |
| 0 expert | 277M | - | - | - | - | - | - |
| 8 experts | - | 191M | 323M | 514M | 195M | 310M | 505M |
| 32 experts | - | 191M | 1,037M | 1,228M | 195M | 1,014M | 1,209M |

# D. Complexity & Parameter Analysis

### D.1. Complexity

We analysis the computational cost of MH-MoE. Without loss of generality, we consider one transformer block with a single-layer SMoE containing $N$ experts, but this can be easily generalized.

The computational cost of SMoE layer comes from two parts: 1) the projection of router. 2) linear projections by parallel experts (FFNs which contains two connected linear layers with inner hidden size $4d$). For a input sequence $\mathbf{X} \in \mathbb{R}^{l \times d}$, the computational cost for each part is

$$\Theta_{\text{router}} = l \times d \times N = ldN, \tag{12}$$

$$\Theta_{\text{experts}} = l \times (d \times 4d + 4d \times d) = 8ld^2. \tag{13}$$

respectively. Thus, the total computational cost of SMoE is $\Theta_{\text{SMoE}} = ld(N + 8d)$.

For MH-MoE layer, in addition to the two computations mentioned above, it also encompasses two additional distinct computational aspects: 3) projection of multi-head layer. 4) projection of merge layer. The computational cost for each part is

$$\Theta_{\text{multi-head}} = l \times d \times d = ld^2, \tag{14}$$

$$\Theta_{\text{router}} = h \times l \times \frac{d}{h} \times N = ldN, \tag{15}$$

$$\Theta_{\text{experts}} = h \times l \times \left( \frac{d}{h} \times 4\beta d + 4\beta d \times \frac{d}{h} \right) = 8\beta ld^2, \tag{16}$$

$$\Theta_{\text{merge}} = l \times d \times d = ld^2, \tag{17}$$

where $\beta$ is a hyperparameter employed to scale the inner hidden dimension of FFNs. In our empirical experiments, we meticulously adjust $\beta$ to ensure a parametric equilibrium between our MH-MoE and SMoE. So, the overall computational cost of MH-MoE is $\Theta_{\text{MH-MoE}} = ld(N + 8\beta d + \frac{N}{h})$.

Thus we compare the computational cost between SMoE and MH-MoE as $\delta$:

$$\delta = \Theta_{\text{SMoE}} - \Theta_{\text{MH-MoE}} = ld \left[ \underbrace{8d(1 - \beta) - \frac{N}{h}}_{\epsilon} \right], \tag{18}$$

In all of our experimental setups, the smallest $\beta$ is $\frac{63}{64}$. Thus $\epsilon$ exists an lower bound when $N = 256$ (set to the largest number of experts), $h = 4$ (set to the fewest head number) and $\beta = \frac{63}{64}$ (set to the minimum $\beta$). In considering this situation, we have $\epsilon = 8 \times 768 \times \frac{1}{64} - \frac{256}{4} = 96 - 64 > 0$. **So we validate that for all experiments setting, we have $\Theta_{\text{SMoE}} - \Theta_{\text{MH-MoE}} > 0$, i.e., the computational cost of our MH-MoE is fewer than SMoE.**

### D.2. Parameter

For SMoE, the parameter contains two parts: 1) the parameter of router. 2) the parameter of parallel experts:

$$\Gamma_{\text{router}} = d \times N, \tag{19}$$

$$\Gamma_{\text{experts}} = N \times (d \times 4d + 4d \times d) = 8d^2 N \tag{20}$$

Thus, the total parameter of SMoE is $\Gamma_{\text{SMoE}} = dN(1 + 8d)$.

For MH-MoE, in addition to the two parts of parameter mentioned above, it also encompasses two additional aspects: 3)

the parameter of multi-head layer. 4) the parameter of merge layer, while the parameter for each part is

$$\Gamma_{\text{head}} = d \times d, \tag{21}$$

$$\Gamma_{\text{router}} = \frac{d}{h} \times N, \tag{22}$$

$$\Gamma_{\text{experts}} = N \times \left( \frac{d}{h} \times 4\beta d + 4\beta d \times \frac{d}{h} \right) = \frac{8\beta d^2 N}{h} \tag{23}$$

$$\Gamma_{\text{merge}} = d \times d, \tag{24}$$

$$\tag{25}$$

So, the overall parameter of MH-MoE is $\Gamma_{\text{MH-MoE}} = 2d^2 + \frac{dN}{h} + \frac{8\beta d^2 N}{h}$. **Detailed parameter comparison can be found in Tabl 11, we ensure that the parameter count of our MH-MoE remains consistent with or lower than that of X-MoE, ensuring a fair and equitable comparison all experiments.**

# E. PyTorch-style Code

We also provide the PyTorch-style code in Algorithm 1 to explain our MH-MoE, which including two main aspects: 1) `Stage 1.` The creation and initialization of multi-head layer and merge layer. 2) `Stage 2.` The segmentation of tokens, followed by processing through an expert network, and ultimately merging.

---

**Algorithm 1** The Overall Procedures of MH-MoE in a PyTorch-like style.

---

**Input:** A MH-MoE model with L parallel SMoE layers M, the number of the experts $k$.

```
# Stage 1: Initial parameter of multi-head layer & merge layer

for i in range(1, L):
    M[i].multi_head_layer = nn.Linear(hidden_dim, hidden_dim)
    M[i].merge_layer = nn.Linear(hidden_dim, hidden_dim)

    # Initialization
    nn.init.xavier_uniform_(M[i].multi_head_layer.weight, gain=1 / math.sqrt(2))
    nn.init.xavier_uniform_(M[i].merge_layer.weight)
    nn.init.constant_(M[i].merge_layer.bias, 0.0)

# Stage 2: The segmentation and merge of tokens for the i-th MH-MoE layer

def MHMoE_Layer(x):
    '''
    Input:
        x : Tensor shape: (batch_size, Length, hidden_dim)
        mask : Tensor shape: (batch_size, Length)

    Output:
        o : Tensor shape: (batch_size, Length, hidden_dim)

    heads: head number of multi_head layer
    '''

    # Processed by multi-head layer
    x = M[i].multi_head_layer(x)

    # Split token & rearrange sub-tokens in parallel
    x = x.reshape(batch_size * Length * heads, hidden_dim // heads).contiguous()
    mask = mask.reshape(-1, 1).repeat(1, heads).reshape(batch_size * Length * heads)

    # Standrad SMoE routing block
    x, mask = router(x, mask)

    # Merge back to the original token form
    x = x.reshape(batch_size * Length, heads, dim // heads).reshape(batch_size * Length,
        dim).contiguous()
    o = M[i].merge_layer(x)

    return o
```

---

*Table 12.* Prompt template for identifying polysemous and false cognates in different languages.

Your role is to identify polysemous and false cognates in different languages from the given textual input (**### Input**). Note that "Polysemous" refers to a word having two or more completely different meanings (for example, "grouse" has meanings related to complaining and also refers to a type of bird), while "false cognates in different languages" refers to words in different languages that have similar forms but carry different meanings (for example, in English and Italian, "camera" looks similar but represents different semantic concepts).

**### Input**
Text: {Prompt}

Note: Please provide your identify results in the following format:

**### Output for Word 1**
Word: [Make sure that there is only a word here.]
Rationale: [Rationale for why this word is polysemous or false cognates, think step by step]

**### Output for Word 2**
Word: [Make sure that there is only a word here.]
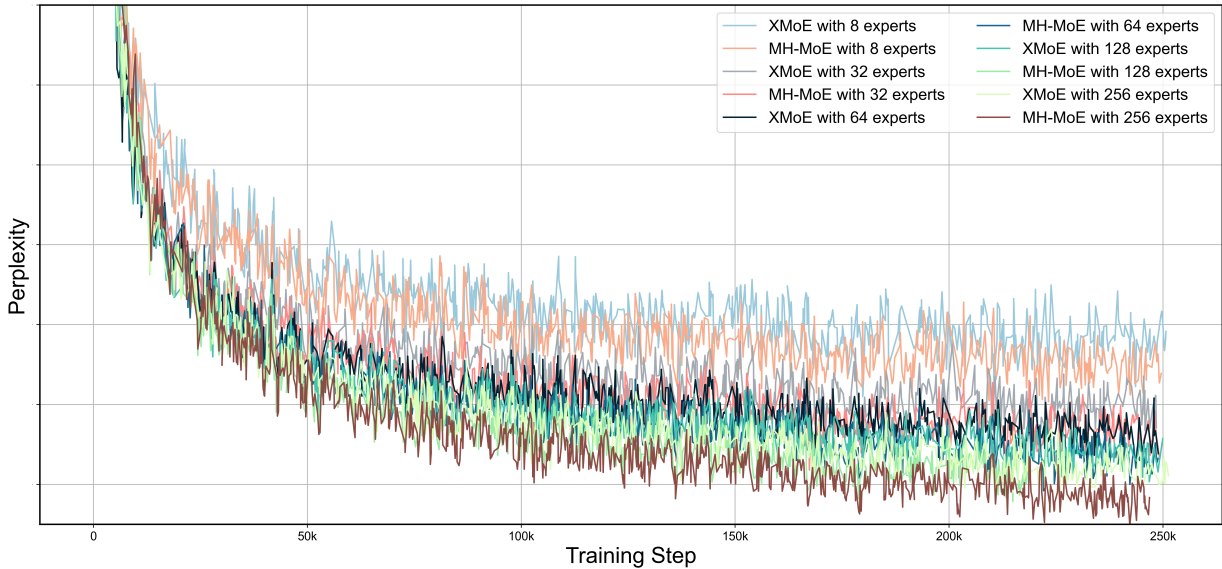Rationale: [Rationale for why this word is polysemous or false cognates, think step by step]
----



*Figure 9.* Validation perplexity reported for both X-MoE and MH-MoE.

# F. Visualization of training perplexity

We provide the training perplexity curve for model training in the experimental setting of increasing the number of experts (from 8 to 256) in Figure 9.

# G. Visualization

We provide more visualization of variation in assign diversity for sub-tokens split from different patches in vision data at Figure 10.
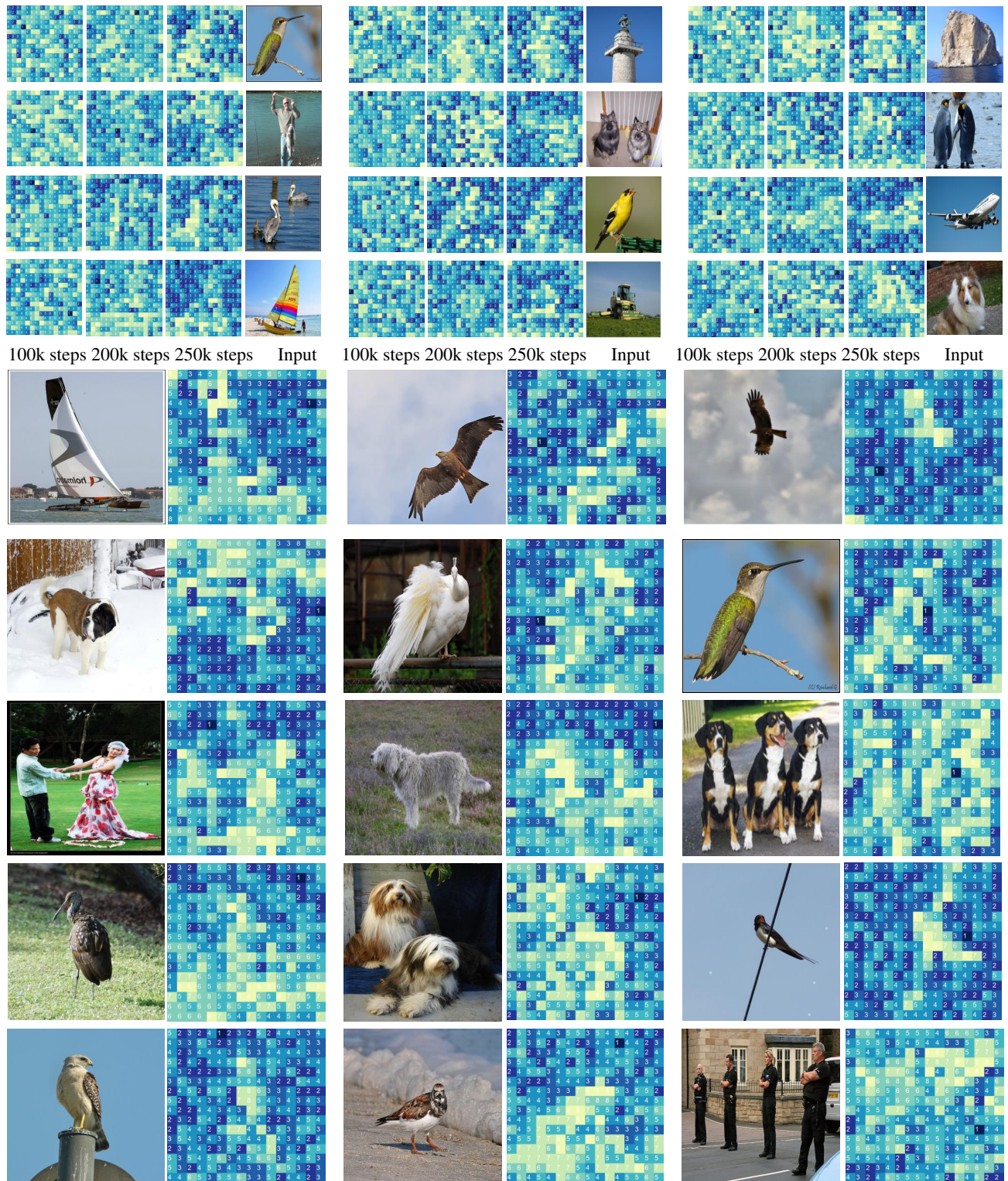
100k steps 200k steps 250k steps     Input     100k steps 200k steps 250k steps     Input     100k steps 200k steps 250k steps     Input

*Figure 10.* More visualization examples for assign diversity of sub-tokens split from different patches with respect to training steps. We analyze MH-MoE with 8 heads ($h$=8) as the subject. Brighter regions indicate that sub-tokens from this patch are distributed to a greater number of diverse experts, while darker regions indicate that sub-tokens are assigned to more of the same experts.