

COMPETESMoE – STATISTICALLY GUARANTEED MIXTURE OF EXPERTS TRAINING VIA COMPETITION

Anonymous authors

Paper under double-blind review

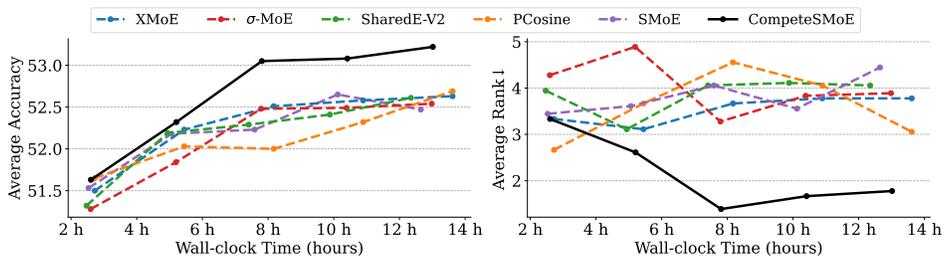


Figure 1: The evolution of zero-shot performance averaged over nine visual instruction tuning tasks throughout training of various SMoE algorithms using a 5.1B parameters backbone.

ABSTRACT

Sparse mixture of experts (SMoE) offers an appealing solution to scale up the model complexity beyond the mean of increasing the network’s depth or width. However, we argue that effective SMoE training remains challenging because of the suboptimal routing process, which often does not involve the experts computation. In this work, we propose *competition*, a novel mechanism to route tokens to experts with the highest neural response. Theoretically, we show that the competition mechanism enjoys a better sample efficiency than the traditional softmax routing. Furthermore, we develop CompeteSMoE, a simple yet effective algorithm for large models by deploying a router to learn the competition policy, thus enjoying strong performances at a low training overhead. Our extensive empirical evaluations on both the visual instruction tuning and language pre-training tasks demonstrate the efficacy, robustness, and scalability of CompeteSMoE compared to state-of-the-art SMoE strategies. We will publish the implementation upon acceptance.

1 INTRODUCTION

Large language models (LLMs) have emerged as a promising architecture for artificial general intelligence. In recent years, LLMs have shown remarkable success in solving many cognitive tasks, ranging from language, vision understanding (Bao et al., 2022b; Gulati et al., 2020; Dosovitskiy et al., 2021; Ruiz et al., 2021; Bao et al., 2022a; Li et al., 2022; 2023a), to code generation (Wang et al., 2021), reinforcement learning (Chow et al., 2023) and life sciences (Rives et al., 2021). Since the release of the original Transformer model (Vaswani et al., 2017), extensive efforts have been devoted to scaling the model complexity to take advantage of massive datasets and advanced computing hardware (Radford et al., 2019; Brown et al., 2020; Du et al., 2022). To go beyond simply increasing the depth and width of the network, Sparse Mixture-of-experts (SMoE) (Fedus et al., 2022) has emerged as an appealing solution for scaling LLMs. By modularizing the network and activating only subsets of experts per input, SMoE offers constant computational costs when increasing the model complexity and often resulting in improved performance.

Despite the initial success, practical SMoE training has been known to be notoriously challenging in both engineering and algorithmic aspects. Thus, despite the rapid development of advanced SMoE research in theory and algorithm (Lee-Thorp & Ainslie, 2022; Riquelme et al., 2021; Chi et al., 2022a), limited progress has been made in leading industrial models such as DeepSeek (DeepSeek-AI et al., 2024; 2025) or Phi-MoE (Abdin et al., 2024) as they still implement variants of the vanilla

054 routing mechanism since the original SMoE (Shazeer et al., 2017; Lepikhin et al., 2021; Fedus et al.,
 055 2022). We argue that this discrepancy exists because many state-of-the-art strategies often rely on
 056 intuitive conceptualizations, which can only offer greedy solutions that work training in the limited
 057 training data and small model regimes. Evidently, many of existing works (Le et al., 2025; Do et al.,
 058 2023; Nguyen et al., 2025; Dai et al., 2022a) still follow the in-domain evaluation and ignore the
 059 zero-shot generalization capabilities of pre-train language models, which are their main use cases.

060 This work makes a step towards a statistically guaranteed SMoE training strategy that can yield
 061 improvements over a wide range of training settings in large-scale models. To this end, we investigate
 062 the core mechanism of routing tokens to experts in SMoE and argue that it could be suboptimal
 063 because the experts performing the calculation do not directly contribute to the routing process.
 064 This limitation has motivated us to develop a radical routing strategy to distribute tokens to experts
 065 more effectively than using the traditional router. To this end, motivated by the Winner-take-all
 066 (WTA) principle (Grossberg & Grossberg, 1982; Riesenhuber & Poggio, 1999; Andersen et al., 1969;
 067 Eccles, 2013), we propose the *competition* mechanism for SMoE training. The core mechanism of
 068 competition is activating all experts and defining a winning criterion so that tokens are only sent to
 069 the winning experts. Thus, competition addresses the fundamental limitation of traditional routing
 070 schemes by involving experts in the routing process, which we rigorously show to achieve a better
 071 sample efficiency or convergence rate than the traditional softmax routing. Furthermore, we go beyond
 072 statistical analysis by developing the *CompeteSMoE* algorithm that implements the competition
 073 mechanism into large-scale models at a modest overhead. Specifically, *CompeteSMoE* improved the
 074 zero-shot performance across 16 common benchmarks in both vision-language finetuning (Figure 1)
 075 and language pre-training settings.

076 In summary, our work makes the following contributions. First, we propose a novel *competition*
 077 mechanism for training SMoE, which enjoys a better convergence rate than softmax routing. Sec-
 078 ond, we develop *CompeteSMoE*, a scalable and effective training strategy for SMoE training via
 079 competition. Lastly, we conduct extensive experiments to explore the behaviours of *CompeteSMoE*,
 080 including its performance, scalability, convergence property, and routing efficacy.

081 2 COMPETESMOE

082 We first recap the foundation of MoE in Section 2.1. Then, we introduce the competition mechanism
 083 in Section 2.2, discuss the scheduled router training in Section 2.3, and detail the *CompeteSMoE*
 084 algorithm in Section 2.4. We provide a list of all notations and their meanings in Table 9, Appendix A.

085 2.1 BACKGROUND

086 The traditional SMoE layer (Shazeer et al., 2017) consists of a router $\mathcal{R}(\cdot, W_r)$ parameterized by W_r
 087 and N experts $\{g(\cdot, W_{e_i})\}_{i=1}^N$ parameterized by $W_{e_i}, i \in [N]$, respectively. The router takes the input
 088 token \mathbf{x} as input and produces an affinity score vector on experts as $\mathbf{s}_{\mathcal{R}} = \sigma(\text{TopK}_{-\infty}(\mathbf{x}^\top W_r))$,
 089 where σ is a scoring function, often implemented as a softmax or sigmoid function. The $\text{TopK}_{-\infty}$
 090 function keeps the largest K elements in a vector and sets the other elements to negative infinity
 091 $(-\infty)$. With this notation, the SMoE layer takes an input token \mathbf{x} and calculate the final output by
 092 aggregating the outputs of each expert weighted by their affinity scores as: $\hat{y} = \sum_{i=1}^N \mathbf{s}_{\mathcal{R}}^i \cdot g(\mathbf{x}; W_{e_i})$
 093 In practice, it is common for K to be smaller than N , i.e. $K < N$, to improve the model efficiency.

094 2.2 ROUTING VIA COMPETITION

095 We now introduce the *competition* mechanism as an effective routing strategy to facilitate SMoE
 096 training. The key idea of competition is allowing all experts to calculate their outputs, and selection
 097 is performed via the winner-take-all mechanism. Thus, experts will compete with one another and
 098 the best ones are selected to calculate the final output. To implement the competition, we propose
 099 to use the expert’s neural response as its affinity score, i.e. $s_i = \mathbb{E}[\kappa(g(\mathbf{x}, W_{e_i}))]$, where $\kappa(\cdot)$ is an
 100 activation function over the expert’s neural responses, and \mathbb{E} denotes the mean over the elements
 101 of the expert’s output vector. In the experiments, we implement κ as the softplus function, unless
 102 otherwise stated. However, our competition mechanism and the theoretical analysis thereafter are
 103 general and do not make strong assumptions about κ . We will provide the results of other choices of
 104 κ in Appendix C.3. With this notation, the training of SMoE with competition is formulated via the
 105 following steps:
 106
 107

1. Compute the output of all N experts for a given input \mathbf{x} as $g(\mathbf{x}, W_{e_i}), \forall i \in [N]$.
2. Compute the affinity score of each expert: $s_i = \mathbb{E}[\log(1 + e^{g(\mathbf{x}, W_{e_i})})], \forall i \in [N]$.
3. Select the Top- K experts based on the highest neural response and compute the normalized affinity scores: $\hat{s}_C^i = \text{TopK}_0(s_i, K), s_C^i = \frac{\hat{s}_C^i}{\sum_{j=1}^N \hat{s}_C^j}$. Here, TopK_0 is similar to the traditional $\text{TopK}_{-\infty}$ but sets the values outside the K highest values to be 0 instead of $-\infty$.
4. Compute the final output as a weighted sum of the selected experts:

$$\hat{y} = \sum_{i=1}^N s_C^i \cdot g(\mathbf{x}, W_{e_i}).$$

Competition starkly contrasts with the standard SMoE implementation discussed in Section 2.1 where the affinity score is calculated as the dot product between the input \mathbf{x} and the columns of the router W_r , then only a few selected experts actually perform their calculation. Although efficient, it results in suboptimal routing policies because the expert selection is detached from the expert’s forward calculation. In contrast, competition proposes that experts who respond the strongest to an input are selected to process that input, while suppressing the other experts. We will rigorously show the theoretical guarantees of routing via competition in Section 3.

2.3 SCHEDULED TRAINING OF THE ROUTER

One drawback of competition-based expert selection is the high computational overhead of activating all experts, which limits its viability to large-scale models. To make competition applicable to LLM training, we propose to incorporate it into the standard router in SMoE. Specifically, we propose a schedule training procedure that periodically trains the router $\mathcal{R}(\cdot; W_r)$ to jointly estimate the competition policy and minimize the task loss. It is important to note that our analysis in Section 3 will show that using the competition policy alone is theoretically sufficient to achieve a faster convergence rate than the vanilla SMoE. In practice, CompeteSmoE in modern architectures stacks many SMoE layers on top of each other, each of which is equipped with a competition mechanism independently. This deep architecture may require significantly more training samples for convergence, which could be much larger than the dataset size and makes training infeasible on our hardware. Therefore, we propose to jointly learn the task loss and match the competition policy to facilitate the router learning. Particularly, without competition activated, the task loss gradient tells the router how to adjust the affinity scores for selected experts only (since inactivated experts do not receive gradients). When competition is active, its gradient tells the router how to adjust the scores for all experts, including those that are not selected to make final predictions. Thus, CompeteSMoE router is expected to facilitate the training and improve the performance. In the following, we present the router loss for effective training and the router schedulers to ensure that training remains efficient.

2.3.1 ROUTER LOSS

The router is trained to learn the competition policy and use it to minimize the task loss. We propose to learn the competition policy by minimizing a distillation loss, $\mathcal{L}_{\mathcal{D}}$, which characterizes the discrepancy between the competition and router policies. For ease of notation, we use $I_C \subset [N]$ to denote the indices of the experts who won the competition. Then, the distillation loss $\mathcal{L}_{\mathcal{D}}$ can be computed by minimizing the mean squared errors (MSE) between the competition and router policies, via their affinity scores as:

$$\mathcal{L}_{\mathcal{D}}(s_{\mathcal{R}}, s_C) = \text{MSE}(s_{\mathcal{R}}, s_C) + \frac{\alpha}{K} \cdot \sum_{j \in I_C} (s_C^j - s_{\mathcal{R}}^j)^2, \quad (1)$$

where $\alpha \in \mathbb{R}^+$ is a hyperparameter to encourage the router to pay more attention to winning experts.

Diversity Loss One of our main experimental settings is using sparse upcycling (Komatsuzaki et al., 2023) to bypass the expensive pre-training cost, which allows us to test SMoE algorithms on larger models with a low budget. However, sparse upcycling duplicates the experts and make them have similar outputs, which results in no competition in the early stages of training and limited training efficacy. To mitigate this issue, we introduce the Diversity Loss, \mathcal{L}_{div} , to promote diverse representations of the winning experts. Formally, given the output matrix $O \in \mathbb{R}^{K \times D}$ representing the outputs of K winning experts for an input \mathbf{x} , the diversity loss is computed as the mean of the

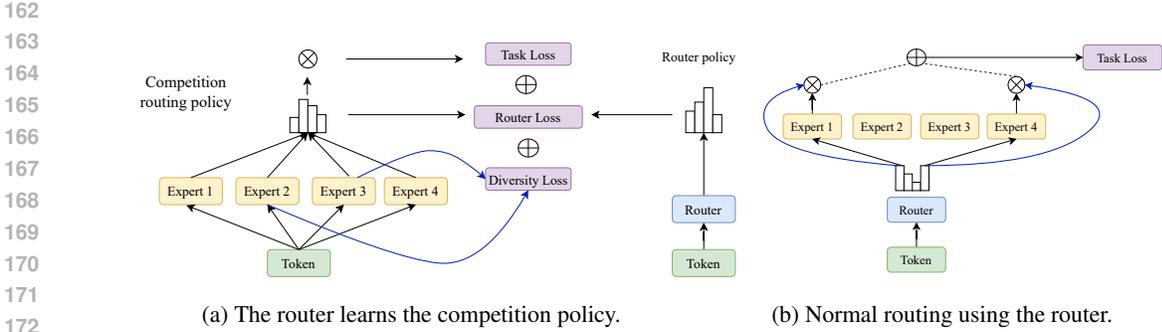


Figure 2: An illustrative of the interleaved learning phases in CompeteSMoE: (a) activating all experts for the router to learn the competition policy; and (b) normal routing using the router.

off-diagonal elements in the correlation matrix constructed from O :

$$\mathcal{L}_{\text{div}}(O) = \frac{1}{K(K-1)} \sum_{i=1}^K \sum_{\substack{j=1 \\ j \neq i}}^K C_{i,j}, \text{ where } C = \frac{O \cdot O^T}{\|O\|_2^2}. \quad (2)$$

We apply the Diversity Loss only within the competition mechanism and emphasize the winning experts as defined in Eq. 2.2, rather than those selected by the router $\mathcal{R}(\cdot; W_r)$. By penalizing winning experts when they produce similar outputs, \mathcal{L}_{div} promotes a more effective competition outcome when using the sparse upcycling strategy.

2.3.2 ROUTER TRAINING SCHEDULE

Schedulers are essential to ensure that the routers can effectively learn a good routing policy while maintaining a limited computational overhead. In the worst case, when all layers of a deep network perform competition simultaneously, this SMoE becomes dense and could crash the training process. Thus, we need to carefully design a schedule to manage the competition frequency across layers. To this end, we employ two schedulers; one is applied per layer independently, while the other monitors the total competition frequency of all layers. For a layer l in a deep network, we first employ a scheduler $\lambda_l(t)$ to determine whether competition should be activated at time step t for this layer. We simply implement $\lambda_l(t)$ by sampling from a Bernoulli distribution with probability ω , which is fixed for all layers. Furthermore, we also employ a global concurrency across layers. Specifically, we only allow the total number of layers performing competition at any time step to be A_{max} . Any layers exceeding this threshold are deferred to perform competition in the next step. Appendix B will provide a detailed formulation of the global scheduler. Based on the number of layers in each model, we set $A_{\text{max}} = 9$ for vision-language models and $A_{\text{max}} = 6$ for the language model pre-training setting, in order to achieve an optimal trade-off between performance and computational feasibility.

2.4 THE COMPETE SMoE ALGORITHM

We are now ready to describe the CompeteSMoE algorithm to enhance SMoE training of large-scale models. Before training, we use the schedulers to generate all time steps for which the competition mechanism is activated at each layer and store them in $\{\Lambda(l)\}_{l=1}^L$, where $\Lambda(l, t) = 1$ indicating that the l -layer will perform competition at time t . Note that this step is performed offline, only one time before training starts. Then, according to the schedule $\Lambda(l, t)$, the training dynamic involves: (i) training the activated experts to minimize the task loss, \mathcal{L}_{NLL} , and (ii) training the activated router to minimize the task and router losses. We provide an illustration of CompeteSMoE training in Figure 2.

We now discuss a general guideline to set the hyper-parameters introduced by CompeteSMoE. We recommend the balancing hyper-parameters α, β, γ to be small values such as 0.01 or 0.005. The Bernoulli parameter ω should also be small (e.g. 0.07) so that competition is not activated too often. The global scheduler threshold should be set based on the specific backbone architecture, model, and training infrastructure to ensure stability. We found $A_{\text{max}} = 9$ for vision-language models and $A_{\text{max}} = 6$ for language model pre-training to maximize the memory usage of our hardware. Lastly, we emphasize that the value ranges of these hyper-parameters can be derived by their definition,

which greatly reduces the effort for hyper-parameter searching. As long as they follow this guideline, we empirically validate the effectiveness of these guidelines through an extensive ablation study in Appendix C, showing that they consistently lead to strong and stable performance in all settings.

3 STATISTICAL GUARANTEE OF THE COMPETITION MECHANISM

In this section, we perform a convergence analysis of Gaussian MoE models equipped with the competition mechanism. Our primary objective is to theoretically justify the effectiveness of the competition mechanism by investigating its sample efficiency in terms of expert estimation.

Problem setting. Let $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n) \in \mathcal{X} \times \mathcal{Y}$ be i.i.d samples drawn from bounded subsets $\mathcal{X} \subset \mathbb{R}^{d_1}$ and $\mathcal{Y} \subset \mathbb{R}$ according to the following conditional density function:

$$p_{G_*}(Y|X) := \sum_{i=1}^{N^*} \frac{\exp(\log(1 + \exp(g(X, W_{e_i}^*)))}{\sum_{j=1}^{N^*} \exp(\log(1 + \exp(g(X, W_{e_j}^*)))} \cdot f(Y|g(X, W_{e_i}^*), \nu_i^*). \quad (3)$$

Here, N^* is the number of ground-truth experts denoted by $g(X, W_{e_i}^*)$, while $f(\cdot|\mu, \nu)$ stands for the Gaussian density with mean μ and variance ν . In addition, we also define $G_* := \sum_{i=1}^{N^*} \delta_{(W_{e_i}^*, \nu_i^*)}$ as a mixing measure with ground-truth parameters $(W_{e_i}^*, \nu_i^*)$, where δ denotes the Dirac measure. For the sake of theory, we assume that $(W_{e_1}^*, \nu_1^*), (W_{e_2}^*, \nu_2^*), \dots, (W_{e_{N^*}}^*, \nu_{N^*}^*)$ are distinct parameters belonging to a compact space $\Theta \subset \mathbb{R}^{d_2} \times \mathbb{R}_+$ for some $d_2 \in \mathbb{N}$. Next, we assume that the expert function $g(X, W_e)$ is non-zero and differentiable with respect to its parameter W_e for almost surely X . Furthermore, for any parameter $W_e \in \mathbb{R}^{d_2}$, if there exists $\alpha_1^{(u)}, \alpha_2^{(uv)}, \alpha_3^{(uv)} \in \mathbb{R}$ for $1 \leq u, v \leq d_2$ such that $\sum_{u=1}^{d_2} \alpha_1^{(u)} \frac{\partial g}{\partial W_e^{(u)}}(X, W_e) + \sum_{u,v=1}^{d_2} \alpha_2^{(uv)} \frac{\partial^2 g}{\partial W_e^{(u)} \partial W_e^{(v)}}(X, W_e) + \sum_{u,v=1}^{d_2} \alpha_3^{(uv)} \frac{\partial g}{\partial W_e^{(u)}}(X, W_e) \frac{\partial g}{\partial W_e^{(v)}}(X, W_e) = 0$ for almost surely X , then we must have $\alpha_1^{(u)} = \alpha_2^{(uv)} = \alpha_3^{(uv)} = 0$ for all $1 \leq u, v \leq d_2$. For example, it can be verified that feed-forward networks (FFNs) of the form $g(X, (W_{e,2}, W_{e,1}, b)) = W_{e,2} \text{Softplus}(W_{e,1}^\top X + b)$ we used in Section 2.2 satisfy this algebraic independence condition. On the other hand, since linear experts $g(X, (a, b)) = a^\top X + b$ does not meet this condition, we will conduct a separate convergence analysis for them in Appendix J.

Maximum likelihood estimation. Since the number of ground-truth experts N^* is typically unknown in practice, we fit the model equation (3) with a mixture of $N > N^*$ experts. Then, we estimate the unknown parameters $(W_{e_i}^*, \nu_i^*)$, for $1 \leq i \leq N$, via estimating the ground-truth mixing measure G_* using the maximum likelihood method as follows:

$$\hat{G}_n \in \arg \max_{G \in \mathcal{G}_N(\Theta)} \frac{1}{n} \sum_{i=1}^n \log(p_G(Y_i|X_i)), \quad (4)$$

where we define $\mathcal{G}_N(\Theta) := \{G = \sum_{i=1}^{N'} \delta_{(W_{e_i}, \nu_i)} : 1 \leq N' \leq N, (W_{e_i}, \nu_i) \in \Theta\}$.

Proposition 3.1. *With the MLE defined in equation (4), the convergence rate of the density estimation $p_{\hat{G}_n}(Y|X)$ to the ground-truth density $p_{G_*}(Y|X)$ is given by:*

$$\mathbb{E}_X[V(p_{\hat{G}_n}(\cdot|X), p_{G_*}(\cdot|X))] = \mathcal{O}_P(\sqrt{\log(n)/n}),$$

Above, we denote $V(p_1, p_2) := \frac{1}{2} \int |p_1 - p_2| dm$ as the Total Variation distance between two probability density functions p_1, p_2 dominated by the Lebesgue measure m .

The proof of Proposition 3.1 can be found in Appendix K.3. The above result indicates that the density estimation $p_{\hat{G}_n}$ converges to its true counterpart p_{G_*} at a parametric rate of order $\tilde{\mathcal{O}}_P(n^{-1/2})$.

Thus, if we can construct some loss function between two mixing measures \hat{G}_n and G_* , denoted by $\mathcal{L}(\hat{G}_n, G_*)$, such that $\mathbb{E}_X[V(p_{\hat{G}_n}(\cdot|X), p_{G_*}(\cdot|X))] \gtrsim \mathcal{L}(\hat{G}_n, G_*)$, then we will obtain parameter and expert estimation rates via the bound $\mathcal{L}(\hat{G}_n, G_*) = \mathcal{O}_P(\sqrt{\log(n)/n})$. For that purpose, let us introduce the concept of Voronoi loss proposed in Manole et al. Manole & Ho (2022).

Voronoi loss. For an arbitrary mixing measure G , we distribute its atoms to the following Voronoi cells generated by the support points of the ground-truth mixing measure G_* :

$$\mathcal{C}_j \equiv \mathcal{C}_j(G) := \{i \in [N] : \|\theta_i - \theta_j^*\| \leq \|\theta_i - \theta_\ell^*\|, \forall \ell \neq j\}, \quad (5)$$

where we denote $\theta_i := (W_{e_i}, \nu_i)$ and $\theta_j^* := (W_{e_j}^*, \nu_j^*)$ for all $i \in [N]$ and $j \in [N^*]$. Here, the cardinality of each Voronoi cell \mathcal{C}_j indicates the number of fitted atoms for the ground-truth atom θ_j^* . Then, we build a loss function based on these Voronoi cells as follows:

$$\begin{aligned} \mathcal{L}_1(G, G_*) := & \sum_{j=1}^{N^*} \left| \sum_{i \in \mathcal{C}_j} \exp(c_i) - \exp(c_j^*) \right| + \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i) \left[\|W_{e_i} - W_{e_j}^*\| + |\nu_i - \nu_j^*| \right] \\ & + \sum_{j \in [N^*]: |\mathcal{C}_j| > 1} \sum_{i \in \mathcal{C}_j} \exp(c_i) \left[\|W_{e_i} - W_{e_j}^*\|^2 + |\nu_i - \nu_j^*|^2 \right]. \quad (6) \end{aligned}$$

Given the above Voronoi loss, we are ready to capture the convergence rates of parameter estimation and expert estimation in Theorem 3.2 whose proof can be found in Appendix K.1.

Theorem 3.2. *The following lower bound holds for any mixing measure $G \in \mathcal{G}_N(\Theta)$:*

$$\mathbb{E}_X[V(p_G(\cdot|X), p_{G_*}(\cdot|X))] \gtrsim \mathcal{L}_1(G, G_*). \quad (7)$$

This lower bound and the result of Theorem 3.1 imply that $\mathcal{L}_1(\widehat{G}_n, G_) = \mathcal{O}_P(\sqrt{\log(n)/n})$.*

A few remarks regarding Theorem 3.2 are in order.

(i) *Expert estimation rates.* From the above results and the formulation of the Voronoi loss \mathcal{L}_1 , it follows that the rates for estimating exact-specified parameters $W_{e_j}^*, \nu_j^*$, i.e., for $j \in [N^*] : |\mathcal{C}_j| = 1$, are of parametric order $\widetilde{\mathcal{O}}_P(n^{-1/2})$. Meanwhile, those for over-specified parameters $W_{e_j}^*, \nu_j^*$, i.e., for $j \in [N^*] : |\mathcal{C}_j| > 1$, are slightly slower, of order $\widetilde{\mathcal{O}}_P(n^{-1/4})$. Since the expert function $g(X, W_e)$ is Lipschitz continuous w.r.t its parameter W_e , we have $|g(X, \widehat{W}_{e_i}^n) - g(X, W_{e_j}^*)| \lesssim \|\widehat{W}_{e_i}^n - W_{e_j}^*\|$ for almost surely X . As a result, the estimation rates for exact-specified and over-specified experts $g(X, W_{e_j}^*)$ are also of orders $\widetilde{\mathcal{O}}_P(n^{-1/2})$ and $\widetilde{\mathcal{O}}_P(n^{-1/4})$, respectively. Furthermore, we show in Appendix J that experts of linear form $g(X, (a, b)) = a^\top X + b$ also admit these estimation rates.

(ii) *Sample efficiency of the competition mechanism.* Therefore, we need at most $\mathcal{O}(\epsilon^{-4})$ data points to approximate these experts with a given error $\epsilon > 0$. On the other hand, when not using the competition mechanism Nguyen et al. (2023a), the convergence rates of expert estimation become significantly slow and decrease when the number of fitted experts increases. For instance, if an expert $g(X, W_{e_j}^*)$ is fitted by three experts, i.e., $|\mathcal{C}_j| = 3$, then its estimation rate is of order $\widetilde{\mathcal{O}}_P(n^{-1/12})$. Thus, we need much more data points, specifically $\mathcal{O}(\epsilon^{-12})$, to approximate this expert. Consequently, we conclude that the competition mechanism improves the sample efficiency in terms of expert estimation.

4 RELATED WORK

Mixture of Experts (MoE) is a fundamental model in machine learning (Jacobs et al., 1991; Jordan & Jacobs, 1994) and an instance of the conditional computation framework where different experts are responsible for different regions of the input space (Yuksel et al., 2012; Bengio, 2013; Masoudnia & Ebrahimpour, 2014; Nguyen & Chamroukhi, 2018; Nguyen, 2021). Extensive efforts have been devoted to establishing a theoretical foundation for MoE, including the universal approximation properties (Norets, 2010; Nguyen et al., 2016; 2019; 2020; 2021a; 2023b), model selection criterion (Khalili, 2010; Montuelle & Le Pennec, 2014; Nguyen et al., 2021b; 2022; 2023c), convergence rate for density estimations (Mendes & Jiang, 2012; Norets & Pelenis, 2021; 2022) and the problem of parameter estimation (Ho et al., 2022; Nguyen et al., 2023a; 2024b;a). SMoE, the sparse variant of MoE, is more commonly applied to scale large language models (Fedus et al., 2022). It is often the architecture of choice in many leading industrial models such as Mixtral (Jiang et al., 2024) and DeepSeek (Dai et al., 2024; DeepSeek-AI et al., 2024; 2025). Within the research community, developing novel routing strategies has been a major focus. Notable strategies include letting experts select tokens (Zhou et al., 2022), improving the expert selection process (Lepikhin et al., 2021; Fedus et al., 2022; Zuo et al., 2022; Chi et al., 2022a; Dai et al., 2022b; Chen et al., 2023; Do et al., 2023), or a global expert assignment scheme (Lewis et al., 2021; Clark et al., 2022). Despite the promising progress, many such strategies often do not scale well to LLMs with billions of parameters or the language pre-training setting. In contrast, our work goes beyond both the pure theoretical or analytical studies by developing a theoretically-grounded algorithm for effective training of large-scale LLM

models inspired by competitive learning, which models neural systems where only the most responsive units activate while suppressing others (McClelland et al., 1987). This principle has historically driven progress across diverse learning paradigms, including self-organization (Von der Malsburg, 1973; Kohonen, 1982), feature discovery (Rumelhart & Zipser, 1985), spiking models (Oster & Liu, 2005), and modern competitive architectures such as maxout networks (Goodfellow et al., 2013) and compete-to-compute (Srivastava et al., 2013). Orthogonal to the aforementioned papers, GShard (Lepikhin et al., 2021) developed an efficient framework to automatically sharding massive SMoE models across many devices. Lastly, sparse upcycling (Komatsuzaki et al., 2023) duplicated pre-trained models to build an MoE, which bypasses the expensive costs of training from scratch.

5 EXPERIMENT

5.1 EXPERIMENTAL SETTINGS

Tasks. We evaluate all methods on two challenging tasks: (i) visual instruction tuning (VIT) and (ii) language model pretraining. For VIT, we adopt the CuMo (Li et al., 2024) and LibMoE framework (Nguyen et al., 2024c), which follows a three-stage training pipeline: pre-training (PT), pre-finetuning (PFT), and visual instruction tuning (VIT). The first two stages are trained with a dense model. In the final VIT stage, sparse upcycling (Komatsuzaki et al., 2023) is applied by resuming from the PFT checkpoint and replacing selected MLP layers in the vision encoder and connector with SMoE blocks. Since one can easily replace the LLMs by other MoE models, the SMoE components are only the vision encoder and the vision-language connector. For the language pretraining task, we adopt the MoEUT (Csordás et al., 2024) framework under the **large setting** and train the SMoE models from scratch. While pre-training has been commonly explored for benchmarking SMoE algorithms (Csordás et al., 2024), it is expensive to scale to large models. Thus, we include the VIT task, which is an emerging and challenging setting that take advantage of pre-training checkpoints, allowing us to evaluate SMoE at a modest cost.

Training. For VIT, we follow Li et al. (2024) to use the LLaVA-558K (Liu et al., 2023a) for PT, ALLaVA (Chen et al., 2024a) for PFT, and LLaVA-665K (Liu et al., 2024a) for VIT. The total tokens for all stages is over 1B. We use Phi-3.5 Mini (Abdin et al., 2024) as the language model and SigLIP (Zhai et al., 2023) as the vision encoder, totaling 5.1B parameters. All MoE algorithms are applied during the VIT stage. We set $N = 4$ experts per layer and activate $K = 2$ experts per token. Training default uses both the balancing and z-losses (Fedus et al., 2022). For language pre-training, we follow MoEUT (Csordás et al., 2024) and use 13B tokens from the SlimPajama corpus (Soboleva et al., 2023). We implement a 1B-parameter decoder-only model, where each SMoE layer contains 24 experts with $K = 8$ experts active per token, and the balancing loss (Fedus et al., 2022). All experiments are conducted on $4 \times \text{H100}$ GPUs with a fixed random seed.

Evaluation Benchmarks. All models are evaluated under the zero-shot settings using the well-established benchmarks from the community. For the VIT task, we consider the following benchmarks: AI2D (Kembhavi et al., 2016), TextVQA Validation (Singh et al., 2019), GQA (Hudson, 2019), HallusionBench (Guan et al., 2023), MathVista (test-mini split) (Lu et al., 2023), MMBench (English subset, dev version) (Liu et al., 2023b), MME RealWorld Lite (Zhang et al., 2025b), MMMU Validation (Yue et al., 2023), MMStar (Chen et al., 2024b), POPE (Li et al., 2023b), and OCR-Bench (Liu et al., 2024b). For benchmarks requiring GPT-based evaluation, such as MathVista and HallusionBench, we use GPT-4o (2024-08-06). These benchmarks are selected to cover a wide range of capabilities of the model, from perception, reasoning, to assessing hallucination. For the language pretraining task, we evaluate on LAMBADA (Paperno et al., 2016), BLiMP (Warstadt et al., 2023), Children’s Book Test (CBT) (Zhang et al., 2025a), HellaSwag (Zellers et al., 2019), PIQA (Bisk et al., 2019), ARC-Easy (Clark et al., 2018), RACE (Lai et al., 2017), and SIQA (Sap et al., 2019), which are commonly used for models at this scale.

Baseline. We compare CompeteSMoE against a suite of state-of-the-art SMoE algorithms. First, SMoE (Fedus et al., 2022), the original SMoE and still stands strong in today’s leading models. Then, we consider activation-based SMoE such as XMoe (Chi et al., 2022b), Perturbed Cosine Router (PCosine) (Nguyen et al., 2025), and σ -MoE (Csordás et al., 2023), which incorporate cosine similarity or sigmoid activation to improve routing efficiency. Furthermore, inspired by the DeepSeek V2 architecture (DeepSeek-AI et al., 2024), we also considered the SharedExpert V2 (SharedE-V2) baseline, which enhances SMoE with one shared expert. Similarly, for the language pretraining task, we also implement the SharedE-V3 baseline, which follows the DeepSeek V3 architecture (DeepSeek-AI et al., 2025). SharedE-V3 replaces the softmax routing in SharedE-V2

with the normalized sigmoid. We use the same hyper-parameter configuration as described above to validate the effectiveness of different SMoE algorithms.

5.2 MAIN RESULTS

Table 1: Performance comparison of SMoE strategies in the ViT sparse upcycling setting with a 5B-parameter model. **Bold** values denote the best results, while underlined values indicate the second best. Symbols \uparrow and \downarrow indicate that higher or lower values are better, respectively.

Method	AI2D	Text VQA	GQA	MM Bench	Hallusion	Math Vista	MMMU	MMStar	POPE	OCR	MME RWL	Avg. Acc \uparrow	Avg. Rank \downarrow
SMoE (Fedus et al., 2022)	65.90	41.23	60.96	70.88	39.64	31.40	42.22	40.52	86.56	32.10	31.89	49.39	4.55
XMoE (Chi et al., 2022a)	65.19	41.14	60.63	71.31	41.22	31.50	42.89	42.60	86.12	31.30	<u>32.51</u>	<u>49.67</u>	3.50
PCosine (Nguyen et al., 2025)	65.45	<u>41.68</u>	<u>61.38</u>	<u>71.56</u>	40.27	30.80	<u>42.56</u>	41.87	86.90	30.80	32.05	49.57	<u>3.42</u>
σ -MoE (Csordás et al., 2023)	65.09	41.37	61.48	71.39	<u>41.01</u>	31.90	41.78	42.10	86.52	32.20	30.95	49.62	3.64
SharedE-V2 (DeepSeek-AI et al., 2024)	64.93	41.53	61.15	71.05	<u>41.20</u>	31.20	<u>42.56</u>	41.44	86.08	<u>32.40</u>	32.36	49.63	4.05
CompeteSMoE	66.22	41.92	61.25	72.59	41.22	<u>31.70</u>	42.00	<u>42.25</u>	86.91	33.20	32.52	50.16	1.77

Table 2: Performance comparison of SMoE strategies in the language pretraining setting with a 1B-parameter model. **Bold** values denote the best results, underlined values indicate the second best.

Method	LAMBADA	BLiMP	CBT	HellaSwag	PIQA	ARC-E	RACE	SIQA	Avg. Acc \uparrow	Avg. Rank \downarrow
SMoE (Fedus et al., 2022)	41.24	80.68	90.63	39.17	65.18	39.66	34.53	38.28	53.67	3.81
XMoE (Chi et al., 2022a)	42.23	80.40	90.44	38.63	64.04	38.60	34.26	37.31	53.24	5.88
PCosine (Nguyen et al., 2025)	41.90	80.35	90.26	38.70	63.71	39.66	34.29	38.13	53.38	5.31
σ -MoE (Csordás et al., 2023)	<u>42.39</u>	80.64	90.63	39.12	64.96	39.66	33.81	<u>38.33</u>	53.69	3.88
SharedE-V2 (DeepSeek-AI et al., 2024)	41.65	80.65	90.63	39.57	65.73	39.15	34.71	37.89	53.75	3.75
SharedE-V3 (DeepSeek-AI et al., 2025)	41.91	80.23	91.02	39.19	65.45	<u>39.53</u>	<u>34.86</u>	37.97	<u>53.77</u>	<u>3.63</u>
CompeteSMoE	42.66	80.92	<u>90.91</u>	<u>39.35</u>	65.91	39.20	34.91	38.43	54.04	1.75

5.2.1 PERFORMANCE COMPARISON

We report the results of the ViT and language pre-training settings in Table 1 and Table 2, respectively. In general, we observe that CompeteSMoE offers significant improvements over many benchmarks. In addition, CompeteSMoE demonstrated the best performance in many of the challenging and important capabilities such as real-world visual perception and reasoning (MME RWL), reducing visual hallucination (Hallusion, POPE), OCR (OCRBench) and commonsense reasoning (PIQA, SIQA). Furthermore, we report the evolution of zero-shot performances of ViT benchmarks throughout training in Figure 1. The results showed that CompeteSMoE consistently achieved better results than the baselines throughout training. Notably, CompeteSMoE demonstrated a significant improvement in training efficiency, where the checkpoint at eight hours (8h) already outperformed all baselines at their final checkpoint of 14 hours. Lastly, we emphasize that the improvements of CompeteSMoE can be considered significant in the zero-shot evaluation setting because its power law indicates that reducing (zero-shot) errors requires a substantial increase in data and compute (Hoffmann et al., 2022; Cherti et al., 2023). Since we fixed the training data, the zero-shot improvements observed purely came from the advanced CompeteSMoE algorithm. Overall, the results corroborate our theoretical results that CompeteSMoE achieved a better sample efficiency and better zero-shot generalization.

5.2.2 EXPERT ROUTING BEHAVIOR ANALYSIS

Table 3: Performance of SMoE and CompeteSMoE when changing top-1 expert to top-(K+1). Numbers in parentheses indicates the changes compared to the original routing results in Table 1.

Method	Text VQA	MMBench	MMMU	MMStar	POPE	OCR Bench	Avg. Change
SMoE	41.09 (-0.14)	71.39 (+0.52)	43.22 (+1.00)	42.94 (+2.42)	86.40 (-0.16)	31.50 (-0.60)	0.51
CompeteSMoE	41.48 (-0.45)	71.22 (-1.37)	41.67 (-0.33)	40.55 (-1.70)	86.10 (-0.81)	31.70 (-1.50)	-1.03

(a) **Evaluating the Effectiveness of Expert Routing.** We investigate the experts selection’s quality of different policies. To this end, during inference, we replace the expert with the highest affinity score with the expert with the $K + 1$ highest score, which is equivalent to shifting the selected experts down by one rank. Table 3 reports the results of this experiment in the ViT setting. The results show that the SMoE routing policy is clearly suboptimal since selecting a worse expert led to improvements on several benchmarks. On the other hand, CompeteSMoE performances drop in all cases when we deliberately deviate from the router that learned the competition policy. This result shows that CompeteSMoE facilitated a more effective routing policy compared to the traditional SMoE.

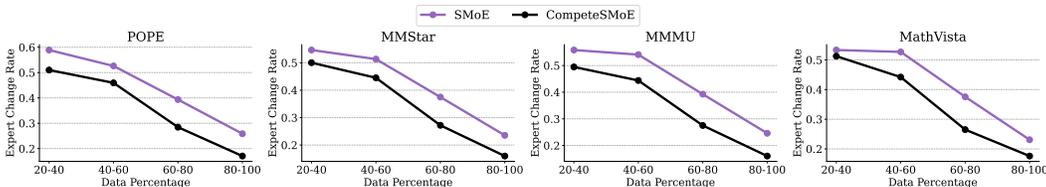


Figure 3: Comparison of expert change rates at different training stages. Lower values are better.

(b) Stability of Expert Routing During Training. We now investigate the router’s convergence rate, showing that CompeteSMoE can quickly find a good routing policy. To this end, we introduce *Expert Change Rate (ECR)* to measure the convergence rate of routers. Specifically, given a dataset \mathcal{D} , we record the expert assignments in all layers for each token in \mathcal{D} using two model checkpoints at time steps T and T' . Then, the ECR of \mathcal{D} from T to T' is the number of mismatched assignments normalized by all assignments. We expect ECR at convergence to be low while high ECR values indicate that the router’s policy is changing and unstable. Figure 3 reports the ECR throughout training on four VIT zero-shot benchmarks. We can see that CompeteSMoE has a lower ECR in all cases, indicating that its routers have a faster convergence rate. This results further support the faster convergence rate and better performance of CompeteSMoE observed in Figure 1 and Table 1.

5.3 COMPLEXITY ANALYSIS

We compare the computational complexities of various methods in Table 4. We report the wall-clock training time, training throughput, inference throughput, and peak GPU memory (excluding cached memory blocks) in the VIT setting of the 5.1B model. The results show that CompeteSMoE’s training overhead compared to SMoE is almost negligible. During inference, CompeteSMoE only uses the simple router, which is exactly the same as SMoE, and is more efficient than cosine similarity-based strategies such as XMoE and PCosine because they introduce additional parameters to the router. CompeteSMoE also incurs a slightly higher peak memory usage compared to other baselines (up to +5%), which was affordable by our hardware. In general, users can adjust the competition’s concurrency threshold A_{max} to achieve a good trade-off between efficiency and efficacy. In summary, this result shows that CompeteSMoE can effectively leverage competition to improve the result at a modest overhead.

Table 4: Computation complexities of various SMoE algorithms.

Method	Training	Throughput (<i>samples/s</i>)		Peak Mem (GB)
	Time	Train	Infer	
SMoE	12h39m	14.59	9.87	43.86
XMoE	13h37m	13.57	8.97	44.02
σ -MoE	12h59m	14.23	9.61	43.93
PCosine	13h37m	13.57	8.59	44.12
SharedE-V2	12h21m	14.95	9.66	42.29
CompeteSMoE	13h01m	14.18	9.88	46.45

5.4 HYPERPARAMETER SENSITIVITY ANALYSIS

We conduct a thorough hyperparameter sensitivity analysis of CompeteSMoE, focusing exclusively on the most critical hyperparameters ($\omega, \alpha, \gamma, \beta, A_{max}$). All results are reported on 9 vision–language benchmarks using the 5.1B-parameter VLM.

Table 5: Effect of ω on competition training.

ω	Avg. Acc \uparrow	Δ SMoE
3%	52.81	+0.34
5%	52.92	+0.45
7%	53.21	+0.74
9%	52.82	+0.35

Table 6: Effect of α on the distillation loss.

α	Avg. Acc \uparrow	Δ SMoE
0.0	52.92	+0.45
0.1	53.21	+0.74
0.2	52.98	+0.51
0.3	52.87	+0.40

Effect of the Distillation Loss Coefficients. We analyze the influence of the two hyperparameters in the distillation loss: the main scaling factor γ and the auxiliary regularization coefficient α . The

Table 7: Effect of γ on the distillation loss.

γ	Avg. Acc \uparrow	Δ SMoE
0.001	52.74	+0.27
0.01	53.21	+0.74
0.03	52.83	+0.36

Table 8: Effect of β on the diversity loss.

β	Avg. Acc \uparrow	Δ SMoE
0.001	52.79	+0.32
0.005	53.21	+0.74
0.01	52.91	+0.44

coefficient γ controls the overall weight of $\mathcal{L}_{\mathcal{D}}$, while α modulates the strength of the penalty applied to the winning experts. As shown in Table 7 and Table 6, the setting $\gamma = 0.01$ and $\alpha = 0.1$ consistently yields the best results. Notably, $\alpha = 0.1$ provides sufficient regularization to guide the router without dominating the main objective, leading to stable and effective learning.

Analysis of Competition Mechanism Activation Frequency. We next examine how frequently the Competition Mechanism (CM) should be activated during training. Table 5 reports performance under different activation frequencies ω . Small values (e.g., $\omega = 3\%$) underperform, likely due to insufficient competitive pressure. Increasing ω improves performance, with the best accuracy (53.21%) obtained at $\omega = 7\%$. Higher activation rates (e.g., $\omega = 9\%$) do not yield additional gains and may introduce instability, indicating a saturation effect. These results suggest that a moderate activation frequency, with ω typically in the range of 5%–7%, is optimal for balancing competitive learning with stable training dynamics.

Effect of the Diversity Loss Coefficient. The diversity coefficient β regulates the dispersion of expert affinity scores by discouraging overly similar outputs among winning experts. As shown in Table 8, model performance peaks at $\beta = 0.005$, which delivers the largest improvement over the SMoE baseline. Smaller values fail to impose sufficient diversity, whereas larger values over-regularize the router and slightly degrade accuracy. Overall, these results indicate that moderate diversity regularization is optimal, promoting balanced expert utilization without interfering with the primary routing objective.

Key Findings. Across all four hyperparameters ($\omega, \alpha, \gamma, \beta$), CompeteSMoE consistently outperforms the standard SMoE baseline even under suboptimal settings. This robustness demonstrates that the method is not overly sensitive to precise hyperparameter tuning and remains stable across a broad range of values. At the same time, the moderate settings we recommend (activation frequency $\omega=5-7\%$, distillation coefficients $\alpha=0.1$ and $\gamma=0.01$, and diversity coefficient $\beta=0.005$) yield the strongest empirical performance, providing the best balance between competitive pressure, router regularization, and expert diversity. Furthermore, our A_{\max} sensitivity analysis in Appendix G demonstrates that increasing A_{\max} generally leads to more stable and robust model performance, with diminishing returns beyond a moderate threshold. We therefore recommend setting A_{\max} according to the guidelines in Section 2.4, which balance computational efficiency and empirical gains.

6 CONCLUSION

This work proposes competition, a novel strategy to route tokens to experts, and rigorously show that it enjoys a better sample efficiency than softmax routing. Building upon this foundation, we develop CompeteSMoE, an effective algorithm to train large-scale SMoE models with competition at a low computational overhead. Extensive experiments on the visual instruction tuning and language pre-training tasks demonstrate that CompeteSMoE enjoys both a faster convergence rate and final performance on many common zero-shot benchmarks at a minimal overhead.

Despite achieving encouraging results, CompeteSMoE introduces several hyper-parameters, which may increase the cost for hyper-parameter search. In Section 2.4, we provided a guideline for hyper-parameter configuration to alleviate this issue. Algorithmically, CompeteSMoE applies competition on each SMoE layer independently and does not take into account the interactions among experts at different layers. An ideal solution is to perform a graph traversal algorithm through the network depth to determine an optimal expert selection at all layers simultaneously. However, this idea goes beyond the scope of this work, and we will leave it for future studies.

REPRODUCIBILITY STATEMENT

We provide full details of our experimental setup in Section 5 and Appendix H. All necessary code, configuration files are included in the Supplementary Materials. Formal proofs supporting our theoretical claims are presented in Appendix K.

ETHICS STATEMENT

This work focuses on the fundamental research involving a theoretical analysis and a training strategy for Sparse Mixture-of-Experts architectures. Due to the abstract nature of our study, it does not involve human-subject data, privacy-sensitive content, or downstream applications. As such, we do not foresee any issues with respect to fairness, privacy, or societal harm.

REFERENCES

- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Hassan Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Singh Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio Cesar Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Parul Chopra, Allison Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Dan Iter, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Young Jin Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Chen Liang, Weishung Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Xianmin Song, Olatunji Ruwase, Praneetha Vaddamanu, Xin Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Michael Wyatt, Can Xu, Jiahang Xu, Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Cheng-Yuan Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone. *ArXiv*, abs/2404.14219, 2024. URL <https://api.semanticscholar.org/CorpusID:269293048>.
- Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav, et al. Smollm2: When smol goes big—data-centric training of a small language model. *arXiv preprint arXiv:2502.02737*, 2025.
- P Andersen, Gary N Gross, T Lomo, and Ola Sveen. Participation of inhibitory and excitatory interneurons in the control of hippocampal cortical output. In *UCLA forum in medical sciences*, volume 11, pp. 415–465, 1969.
- Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. BEiT: BERT Pre-Training of Image Transformers. In *International Conference on Learning Representations*, 2022a. URL <https://openreview.net/forum?id=p-BhZSz59o4>.
- Hangbo Bao, Wenhui Wang, Li Dong, Qiang Liu, Owais Khan Mohammed, Kriti Aggarwal, Subhojit Som, Songhao Piao, and Furu Wei. VLMO: Unified Vision-Language Pre-Training with Mixture-of-Modality-Experts. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022b. URL <https://openreview.net/forum?id=bydKs84JEyw>.
- Yoshua Bengio. Deep Learning of Representations: Looking Forward. In Adrian-Horia Dediu, Carlos Martín-Vide, Ruslan Mitkov, and Bianca Truthe (eds.), *Statistical Language and Speech Processing*, pp. 1–37, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg. ISBN 978-3-642-39593-2.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language, 2019. URL <https://arxiv.org/abs/1911.11641>.

- 594 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
595 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
596 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
597
- 598 Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang,
599 Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. Allava: Harnessing gpt4v-synthesized
600 data for a lite vision-language model. *arXiv preprint arXiv:2402.11684*, 2024a.
- 601 Lin Chen, Jinsong Li, Xiao wen Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi
602 Wang, Yu Qiao, Dahua Lin, and Feng Zhao. Are we on the right way for evaluating large vision-
603 language models? *ArXiv*, abs/2403.20330, 2024b. URL <https://api.semanticscholar.org/CorpusID:268793433>.
- 604
605
- 606 Tianlong Chen, Zhenyu Zhang, AJAY KUMAR JAISWAL, Shiwei Liu, and Zhangyang Wang.
607 Sparse moe as the new dropout: Scaling dense and self-slimmable transformers. In *The Eleventh
608 International Conference on Learning Representations*, 2023.
- 609
- 610 Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade
611 Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. Reproducible scaling laws for
612 contrastive language-image learning. In *Proceedings of the IEEE/CVF conference on computer
613 vision and pattern recognition*, pp. 2818–2829, 2023.
- 614 Zewen Chi, Li Dong, Shaohan Huang, Damai Dai, Shuming Ma, Barun Patra, Saksham Singhal,
615 Payal Bajaj, Xia Song, Xian-Ling Mao, Heyan Huang, and Furu Wei. On the Representation
616 Collapse of Sparse Mixture of Experts. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave,
617 and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022a. URL
618 <https://openreview.net/forum?id=mWaYC6Czf5>.
- 619
- 620 Zewen Chi, Li Dong, Shaohan Huang, Damai Dai, Shuming Ma, Barun Patra, Saksham Singhal, Payal
621 Bajaj, Xia Song, Xian-Ling Mao, Heyan Huang, and Furu Wei. On the representation collapse of
622 sparse mixture of experts, 2022b. URL <https://arxiv.org/abs/2204.09179>.
- 623 Yinlam Chow, Azamat Tulebergenov, Ofir Nachum, Dhawal Gupta, Moonkyung Ryu, Mohammad
624 Ghavamzadeh, and Craig Boutilier. A Mixture-of-Expert Approach to RL-based Dialogue Man-
625 agement. In *The Eleventh International Conference on Learning Representations*, 2023. URL
626 <https://openreview.net/forum?id=4FBUihxz5nm>.
- 627
- 628 Aidan Clark, Diego De Las Casas, Aurelia Guy, Arthur Mensch, Michela Paganini, Jordan Hoffmann,
629 Bogdan Damoc, Blake Hechtman, Trevor Cai, Sebastian Borgeaud, George Bm Van Den Driess-
630 che, Eliza Rutherford, Tom Hennigan, Matthew J Johnson, Albin Cassirer, Chris Jones, Elena
631 Buchatskaya, David Budden, Laurent Sifre, Simon Osindero, Oriol Vinyals, Marc’Aurelio Ran-
632 zato, Jack Rae, Erich Elsen, Koray Kavukcuoglu, and Karen Simonyan. Unified Scaling Laws for
633 Routed Language Models. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari,
634 Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine
635 Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 4057–4086. PMLR, July
636 2022. URL <https://proceedings.mlr.press/v162/clark22a.html>.
- 637 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
638 Oyvind Taffjord. Think you have solved question answering? try ARC, the AI2 reasoning challenge.
639 *Preprint arXiv:1803.05457*, 2018.
- 640
- 641 Róbert Csordás, Kazuki Irie, and Jürgen Schmidhuber. Approximating two-layer feedforward
642 networks for efficient transformers. *arXiv preprint arXiv:2310.10837*, 2023.
- 643 Róbert Csordás, Kazuki Irie, Jürgen Schmidhuber, Christopher Potts, and Christopher D. Manning.
644 Moeut: Mixture-of-experts universal transformers, 2024. URL <https://arxiv.org/abs/2405.16039>.
- 645
- 646
- 647 Damai Dai, Li Dong, Shuming Ma, Bo Zheng, Zhifang Sui, Baobao Chang, and Furu Wei. Stablemoe:
Stable routing strategy for mixture of experts. *arXiv preprint arXiv:2204.08396*, 2022a.

- 648 Damai Dai, Li Dong, Shuming Ma, Bo Zheng, Zhifang Sui, Baobao Chang, and Furu Wei. StableMoE:
649 Stable Routing Strategy for Mixture of Experts. In *Proceedings of the 60th Annual Meeting of*
650 *the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7085–7095, Dublin,
651 Ireland, May 2022b. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.
652 489. URL <https://aclanthology.org/2022.acl-long.489>.
- 653 Damai Dai, Chengqi Deng, Chenggang Zhao, RX Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding
654 Zeng, Xingkai Yu, Yu Wu, et al. Deepseekmoe: Towards ultimate expert specialization in mixture-
655 of-experts language models. *arXiv preprint arXiv:2401.06066*, 2024.
- 656
657 DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi
658 Deng, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li,
659 Fangyun Lin, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Hanwei Xu, Hao
660 Yang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian
661 Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jin Chen, Jingyang Yuan, Junjie Qiu, Junxiao Song, Kai
662 Dong, Kaige Gao, Kang Guan, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Liyue
663 Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming
664 Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J.
665 Chen, R. L. Jin, Ruiqi Ge, Ruizhe Pan, Runxin Xu, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan
666 Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou,
667 Shuiping Yu, Shunfeng Zhou, Size Zheng, T. Wang, Tian Pei, Tian Yuan, Tianyu Sun, W. L.
668 Xiao, Wangding Zeng, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wentao Zhang, X. Q. Li,
669 Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang
670 Chen, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Liu, Xin Xie, Xingkai
671 Yu, Xinnan Song, Xinyi Zhou, Xinyu Yang, Xuan Lu, Xuecheng Su, Y. Wu, Y. K. Li, Y. X. Wei,
672 Y. X. Zhu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui
673 Wang, Yi Zheng, Yichao Zhang, Yiliang Xiong, Yilong Zhao, Ying He, Ying Tang, Yishi Piao,
674 Yixin Dong, Yixuan Tan, Yiyuan Liu, Yongji Wang, Yongqiang Guo, Yuchen Zhu, Yudian Wang,
675 Yuheng Zou, Yukun Zha, Yunxian Ma, Yuting Yan, Yuxiang You, Yuxuan Liu, Z. Z. Ren, Zehui
676 Ren, Zhangli Sha, Zhe Fu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhewen Hao, Zhihong Shao,
677 Zhiniu Wen, Zhipeng Xu, Zhongyu Zhang, Zhuoshu Li, Zihan Wang, Zihui Gu, Zilin Li, and Ziwei
678 Xie. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model, 2024.
679 URL <https://arxiv.org/abs/2405.04434>.
- 680 DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang
681 Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli
682 Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen,
683 Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding,
684 Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi
685 Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song,
686 Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang,
687 Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan
688 Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang,
689 Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi
690 Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li,
691 Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye,
692 Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang,
693 Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanxia Zhao, Wei An, Wen Liu,
694 Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang,
695 Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha
696 Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu,
697 Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su,
698 Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong
699 Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng,
700 Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan
701 Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yudian Wang, Yue
Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo,
Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu,
Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou,
Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu

- 702 Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan.
703 Deepseek-v3 technical report, 2025. URL <https://arxiv.org/abs/2412.19437>.
704
- 705 Giang Do, Khiem Le, Quang Pham, Trungtin Nguyen, Thanh-Nam Doan, Bint T Nguyen, Chenghao
706 Liu, Savitha Ramasamy, Xiaoli Li, and Steven Hoi. Hyperrouter: Towards efficient training and
707 inference of sparse mixture of experts. *arXiv preprint arXiv:2312.07035*, 2023.
- 708 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
709 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
710 and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale.
711 In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
712
- 713 Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim
714 Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, Barret Zoph, Liam Fedus, Maarten P Bosma,
715 Zongwei Zhou, Tao Wang, Emma Wang, Kellie Webster, Marie Pellat, Kevin Robinson, Kathleen
716 Meier-Hellstern, Toju Duke, Lucas Dixon, Kun Zhang, Quoc Le, Yonghui Wu, Zhifeng Chen,
717 and Claire Cui. GLaM: Efficient Scaling of Language Models with Mixture-of-Experts. In
718 Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato
719 (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of
720 *Proceedings of Machine Learning Research*, pp. 5547–5569. PMLR, July 2022. URL <https://proceedings.mlr.press/v162/du22c.html>.
721
- 722 John C Eccles. *The cerebellum as a neuronal machine*. Springer Science & Business Media, 2013.
723
- 724 William Fedus, Barret Zoph, and Noam Shazeer. Switch Transformers: Scaling to Trillion Parameter
725 Models with Simple and Efficient Sparsity. *Journal of Machine Learning Research*, 23(120):1–39,
726 2022. URL <http://jmlr.org/papers/v23/21-0998.html>.
727
- 728 Ian Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. Maxout
729 networks. In *International conference on machine learning*, pp. 1319–1327. PMLR, 2013.
730
- 731 Stephen Grossberg and Stephen Grossberg. Contour enhancement, short term memory, and constancies
732 in reverberating neural networks. *Studies of Mind and Brain: Neural Principles of Learning,
733 Perception, Development, Cognition, and Motor Control*, pp. 332–378, 1982.
- 734 Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang
735 Chen, Furong Huang, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. Hallusionbench:
736 An advanced diagnostic suite for entangled language hallucination and visual illusion in large
737 vision-language models. 2023. URL [https://api.semanticscholar.org/CorpusID:
738 265499116](https://api.semanticscholar.org/CorpusID:265499116).
- 739 Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo
740 Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang. Conformer: Convolution-augmented
741 Transformer for Speech Recognition. In *Proc. Interspeech 2020*, pp. 5036–5040, 2020. doi:
742 10.21437/Interspeech.2020-3015.
743
- 744 Nhat Ho, Chiao-Yu Yang, and Michael I. Jordan. Convergence Rates for Gaussian Mixtures of
745 Experts. *Journal of Machine Learning Research*, 23(323):1–81, 2022. URL [http://jmlr.
746 org/papers/v23/20-1129.html](http://jmlr.org/papers/v23/20-1129.html).
- 747 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
748 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al.
749 Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
750
- 751 Drew A. Hudson. Gqa : A new dataset for real-world visual reasoning and compositional question
752 answering – supplementary material. 2019. URL [https://api.semanticscholar.org/
753 CorpusID:268114221](https://api.semanticscholar.org/CorpusID:268114221).
- 754 Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of
755 local experts. *Neural computation*, 3(1):79–87, 1991. Publisher: MIT Press.

- 756 Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris
757 Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand,
758 Gianna Lengyel, Guillaume Bour, Guillaume Lample, L elio Renard Lavaud, Lucile Saulnier, Marie-
759 Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le
760 Scao, Th eophile Gervet, Thibaut Lavril, Thomas Wang, Timoth ee Lacroix, and William El Sayed.
761 Mixtral of experts, 2024. URL <https://arxiv.org/abs/2401.04088>.
- 762 Michael I Jordan and Robert A Jacobs. Hierarchical mixtures of experts and the EM algorithm.
763 *Neural computation*, 6(2):181–214, 1994. Publisher: MIT Press.
- 764 Aniruddha Kembhavi, Michael Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali
765 Farhadi. A diagram is worth a dozen images. *ArXiv*, abs/1603.07396, 2016. URL <https://api.semanticscholar.org/CorpusID:2682274>.
- 766 Abbas Khalili. New estimation and feature selection methods in mixture-of-experts models. *Canadian*
767 *Journal of Statistics*, 38(4):519–539, 2010.
- 768 Teuvo Kohonen. Self-organized formation of topologically correct feature maps. *Biological cyber-*
769 *netics*, 43(1):59–69, 1982.
- 770 Aran Komatsuzaki, Joan Puigcerver, James Lee-Thorp, Carlos Riquelme Ruiz, Basil Mustafa, Joshua
771 Ainslie, Yi Tay, Mostafa Dehghani, and Neil Houlsby. Sparse upcycling: Training mixture-of-
772 experts from dense checkpoints, 2023. URL <https://arxiv.org/abs/2212.05055>.
- 773 Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading
774 comprehension dataset from examinations, 2017.
- 775 Minh Le, An Nguyen, Huy Nguyen, Trang Nguyen, Trang Pham, Linh Van Ngo, and Nhat Ho.
776 Mixture of experts meets prompt-based continual learning, 2025. URL <https://arxiv.org/abs/2405.14124>.
- 777 James Lee-Thorp and Joshua Ainslie. Sparse Mixers: Combining MoE and Mixing to build a more
778 efficient BERT. In *Findings of the Association for Computational Linguistics: EMNLP 2022*,
779 pp. 58–75, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational
780 Linguistics. URL <https://aclanthology.org/2022.findings-emnlp.5>.
- 781 Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang,
782 Maxim Krikun, Noam Shazeer, and Zhifeng Chen. GShard: Scaling Giant Models with Conditional
783 Computation and Automatic Sharding. In *International Conference on Learning Representations*,
784 2021. URL <https://openreview.net/forum?id=qrwe7XHTmYb>.
- 785 Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. BASE
786 Layers: Simplifying Training of Large, Sparse Models. In Marina Meila and Tong Zhang
787 (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139
788 of *Proceedings of Machine Learning Research*, pp. 6265–6274. PMLR, July 2021. URL
789 <https://proceedings.mlr.press/v139/lewis21a.html>.
- 790 Jiachen Li, Xinyao Wang, Sijie Zhu, Chia-Wen Kuo, Lu Xu, Fan Chen, Jitesh Jain, Humphrey Shi,
791 and Longyin Wen. Cumo: Scaling multimodal llm with co-upcycled mixture-of-experts, 2024.
792 URL <https://arxiv.org/abs/2405.05949>.
- 793 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-
794 training for unified vision-language understanding and generation. In *International Conference on*
795 *Machine Learning*, pp. 12888–12900. PMLR, 2022.
- 796 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-
797 training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*,
798 2023a.
- 799 Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji rong Wen. Evaluating object
800 hallucination in large vision-language models. In *Conference on Empirical Methods in Natural*
801 *Language Processing*, 2023b. URL <https://api.semanticscholar.org/CorpusID:258740697>.

- 810 Bin Lin, Zhenyu Tang, Yang Ye, Jiayi Cui, Bin Zhu, Peng Jin, Jinfa Huang, Junwu Zhang, Yatian
811 Pang, Munan Ning, et al. Moe-llava: Mixture of experts for large vision-language models. *arXiv*
812 *preprint arXiv:2401.15947*, 2024.
- 813 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*,
814 2023a.
- 815 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction
816 tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
817 pp. 26296–26306, 2024a.
- 818 Yuezhan Li, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike
819 Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your multi-
820 modal model an all-around player? *ArXiv*, abs/2307.06281, 2023b. URL [https://api.
821 semanticscholar.org/CorpusID:259837088](https://api.semanticscholar.org/CorpusID:259837088).
- 822 Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin,
823 Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large
824 multimodal models. *Science China Information Sciences*, 67(12), December 2024b. ISSN
825 1869-1919. doi: 10.1007/s11432-024-4235-6. URL [http://dx.doi.org/10.1007/
826 s11432-024-4235-6](http://dx.doi.org/10.1007/s11432-024-4235-6).
- 827 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chun yue Li, Hannaneh Hajishirzi, Hao Cheng,
828 Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning
829 of foundation models in visual contexts. In *International Conference on Learning Representations*,
830 2023. URL <https://api.semanticscholar.org/CorpusID:264491155>.
- 831 T. Manole and N. Ho. Refined convergence rates for maximum likelihood estimation under finite
832 mixture models. In *Proceedings of the 39th International Conference on Machine Learning*,
833 volume 162 of *Proceedings of Machine Learning Research*, pp. 14979–15006. PMLR, 17–23 Jul
834 2022.
- 835 Saeed Masoudnia and Reza Ebrahimpour. Mixture of experts: a literature survey. *Artificial Intelli-*
836 *gence Review*, 42(2):275–293, 2014. ISSN 1573-7462. doi: 10.1007/s10462-012-9338-y. URL
837 <https://doi.org/10.1007/s10462-012-9338-y>.
- 838 James L McClelland, David E Rumelhart, PDP Research Group, et al. *Parallel distributed processing*,
839 *volume 2: Explorations in the microstructure of cognition: Psychological and biological models*,
840 volume 2. MIT press, 1987.
- 841 Eduardo F Mendes and Wenxin Jiang. On convergence rates of mixtures of polynomial experts.
842 *Neural computation*, 24(11):3025–3051, 2012. Publisher: MIT Press.
- 843 Lucie Montuelle and Erwan Le Pennec. Mixture of Gaussian regressions model with logistic weights,
844 a penalized maximum likelihood approach. *Electronic Journal of Statistics*, 8(1):1661–1695, 2014.
845 Publisher: The Institute of Mathematical Statistics and the Bernoulli Society.
- 846 Hien D Nguyen and Faicel Chamroukhi. Practical and theoretical aspects of mixture-of-experts
847 modeling: An overview. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*,
848 8(4):e1246, 2018. Publisher: Wiley Online Library.
- 849 Hien D Nguyen, Luke R Lloyd-Jones, and Geoffrey J McLachlan. A universal approximation theorem
850 for mixture-of-experts models. *Neural computation*, 28(12):2585–2593, 2016. Publisher: MIT
851 Press.
- 852 Hien D Nguyen, Faicel Chamroukhi, and Florence Forbes. Approximation results regarding the
853 multiple-output Gaussian gated mixture of linear experts model. *Neurocomputing*, 366:208–214,
854 2019. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2019.08.014>. URL <https://www.sciencedirect.com/science/article/pii/S0925231219311336>.
- 855 Hien Duy Nguyen, TrungTin Nguyen, Faicel Chamroukhi, and Geoffrey John McLachlan. Ap-
856 proximations of conditional probability density functions in Lebesgue spaces via mixture of
857 experts models. *Journal of Statistical Distributions and Applications*, 8(1):13, August 2021a.
858 ISSN 2195-5832. doi: 10.1186/s40488-021-00125-0. URL [https://doi.org/10.1186/
859 s40488-021-00125-0](https://doi.org/10.1186/s40488-021-00125-0).

- 864 Huy Nguyen, TrungTin Nguyen, and Nhat Ho. Demystifying softmax gating in Gaussian mixture of
865 experts. In *Advances in Neural Information Processing Systems*, 2023a.
- 866
- 867 Huy Nguyen, Pedram Akbarian, and Nhat Ho. Is temperature sample efficient for softmax Gaussian
868 mixture of experts? In *Proceedings of the ICML*, 2024a.
- 869 Huy Nguyen, TrungTin Nguyen, Khai Nguyen, and Nhat Ho. Towards convergence rates for parameter
870 estimation in Gaussian-gated mixture of experts. In *Proceedings of The 27th International
871 Conference on Artificial Intelligence and Statistics*, 2024b.
- 872
- 873 Huy Nguyen, Pedram Akbarian, Trang Pham, Trang Nguyen, Shujian Zhang, and Nhat Ho. Statistical
874 advantages of perturbing cosine router in mixture of experts. In *International Conference on
875 Learning Representations*, 2025.
- 876 Nam V. Nguyen, Thong T. Doan, Luong Tran, Van Nguyen, and Quang Pham. Libmoe: A library for
877 comprehensive benchmarking mixture of experts in large language models, 2024c.
- 878
- 879 TrungTin Nguyen. *Model Selection and Approximation in High-dimensional Mixtures of Experts
880 Models: from Theory to Practice*. PhD Thesis, Normandie Université, December 2021. URL
881 <https://tel.archives-ouvertes.fr/tel-03524749>.
- 882 TrungTin Nguyen, Hien D. Nguyen, Faicel Chamroukhi, and Geoffrey J. McLachlan. Approximation
883 by finite mixtures of continuous density functions that vanish at infinity. *Cogent Mathematics &
884 Statistics*, 7(1):1750861, January 2020. ISSN null. doi: 10.1080/25742558.2020.1750861. URL
885 <https://doi.org/10.1080/25742558.2020.1750861>. Publisher: Cogent OA.
- 886 TrungTin Nguyen, Hien D. Nguyen, Faicel Chamroukhi, and Geoffrey J. McLachlan. An l_1 -oracle
887 inequality for the Lasso in mixture-of-experts regression models. *arXiv:2009.10622*, January
888 2021b. URL <http://arxiv.org/abs/2009.10622>.
- 889
- 890 TrungTin Nguyen, Hien Duy Nguyen, Faicel Chamroukhi, and Florence Forbes. A non-asymptotic
891 approach for model selection via penalization in high-dimensional mixture of experts models.
892 *Electronic Journal of Statistics*, 16(2):4742 – 4822, 2022. doi: 10.1214/22-EJS2057. URL
893 <https://doi.org/10.1214/22-EJS2057>.
- 894 TrungTin Nguyen, Faicel Chamroukhi, Hien D. Nguyen, and Geoffrey J. McLachlan. Approximation
895 of probability density functions via location-scale finite mixtures in Lebesgue spaces. *Communica-
896 tions in Statistics - Theory and Methods*, 52(14):5048–5059, 2023b. doi: 10.1080/03610926.2021.
897 2002360. URL <https://doi.org/10.1080/03610926.2021.2002360>.
- 898 TrungTin Nguyen, Dung Ngoc Nguyen, Hien Duy Nguyen, and Faicel Chamroukhi. A non-asymptotic
899 risk bound for model selection in high-dimensional mixture of experts via joint rank and variable
900 selection. In *Australasian Joint Conference on Artificial Intelligence*. Springer, 2023c.
- 901
- 902 Andriy Norets. Approximation of conditional densities by smooth mixtures of regressions. *The
903 Annals of Statistics*, 38(3):1733 – 1766, 2010. doi: 10.1214/09-AOS765. URL [https://doi.
904 org/10.1214/09-AOS765](https://doi.org/10.1214/09-AOS765). Publisher: Institute of Mathematical Statistics.
- 905 Andriy Norets and Justinas Pelenis. Adaptive Bayesian estimation of conditional discrete-continuous
906 distributions with an application to stock market trading activity. *Journal of Econometrics*,
907 2021. ISSN 0304-4076. doi: <https://doi.org/10.1016/j.jeconom.2021.11.004>. URL [https://
908 www.sciencedirect.com/science/article/pii/S030440762100261X](https://www.sciencedirect.com/science/article/pii/S030440762100261X).
- 909 Andriy Norets and Justinas Pelenis. Adaptive Bayesian Estimation of Discrete-Continuous Dis-
910 tributions Under Smoothness and Sparsity. *Econometrica*, 90(3):1355–1377, 2022. doi:
911 <https://doi.org/10.3982/ECTA17884>. URL [https://onlinelibrary.wiley.com/doi/
912 abs/10.3982/ECTA17884](https://onlinelibrary.wiley.com/doi/abs/10.3982/ECTA17884).
- 913 Matthias Oster and Shih-Chii Liu. Spiking inputs to a winner-take-all network. *Advances in Neural
914 Information Processing Systems*, 18, 2005.
- 915
- 916 Bowen Pan, Yikang Shen, Haokun Liu, Mayank Mishra, Gaoyuan Zhang, Aude Oliva, Colin Raffel,
917 and Rameswar Panda. Dense training, sparse inference: Rethinking training of mixture-of-experts
language models, 2024.

- 918 Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi,
919 Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset:
920 Word prediction requiring a broad discourse context, 2016. URL [https://arxiv.org/abs/
921 1606.06031](https://arxiv.org/abs/1606.06031).
- 922 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language
923 models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 924 Maximilian Riesenhuber and Tomaso Poggio. Hierarchical models of object recognition in cortex.
925 *Nature neuroscience*, 2(11):1019–1025, 1999.
- 926 Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André
927 Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts.
928 *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021.
- 929 Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo,
930 Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function
931 emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the
932 National Academy of Sciences*, 118, 2021. URL [https://www.pnas.org/doi/abs/10.
933 1073/pnas.2016239118](https://www.pnas.org/doi/abs/10.1073/pnas.2016239118).
- 934 Carlos Riquelme Ruiz, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton,
935 André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling Vision with Sparse Mixture of
936 Experts. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances in
937 Neural Information Processing Systems*, 2021. URL [https://openreview.net/forum?
938 id=NGPmH3vbAA_](https://openreview.net/forum?id=NGPmH3vbAA_).
- 939 David E Rumelhart and David Zipser. Feature discovery by competitive learning. *Cognitive science*,
940 9(1):75–112, 1985.
- 941 Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialliqa: Commonsense
942 reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- 943 Noam Shazeer, *Azalia Mirhoseini, *Krzysztof Maziarsz, Andy Davis, Quoc Le, Geoffrey Hinton, and
944 Jeff Dean. Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer.
945 In *International Conference on Learning Representations*, 2017. URL [https://openreview.
946 net/forum?id=BlckMDqlg](https://openreview.net/forum?id=BlckMDqlg).
- 947 Fangxun Shu, Yue Liao, Le Zhuo, Chenning Xu, Lei Zhang, Guanghao Zhang, Haonan Shi, Long
948 Chen, Tao Zhong, Wanggui He, et al. Llava-mod: Making llava tiny via moe knowledge distillation.
949 *arXiv preprint arXiv:2408.15881*, 2024.
- 950 Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh,
951 and Marcus Rohrbach. Towards vqa models that can read. *2019 IEEE/CVF Conference on
952 Computer Vision and Pattern Recognition (CVPR)*, pp. 8309–8318, 2019. URL [https://api.
953 semanticscholar.org/CorpusID:85553602](https://api.semanticscholar.org/CorpusID:85553602).
- 954 Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey.
955 SlimPajama: A 627B token cleaned and deduplicated version of RedPajama, June 2023. URL
956 <https://huggingface.co/datasets/cerebras/SlimPajama-627B>.
- 957 Rupesh K Srivastava, Jonathan Masci, Sohrob Kazerounian, Faustino Gomez, and Jürgen Schmidhu-
958 ber. Compete to compute. *Advances in neural information processing systems*, 26, 2013.
- 959 S. van de Geer. *Empirical Processes in M-estimation*. Cambridge University Press, 2000.
- 960 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
961 Lukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In I. Guyon, U. Von
962 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-
963 vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,
964 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/
965 file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).

- 972 Chr Von der Malsburg. Self-organization of orientation sensitive cells in the striate cortex. *Kybernetik*,
973 14(2):85–100, 1973.
- 974
- 975 Yue Wang, Weishi Wang, Shafiq Joty, and Steven CH Hoi. Codet5: Identifier-aware unified pre-trained
976 encoder-decoder models for code understanding and generation. *arXiv preprint arXiv:2109.00859*,
977 2021.
- 978 Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohanane, Wei Peng, Sheng-Fu Wang, and
979 Samuel R. Bowman. Blimp: The benchmark of linguistic minimal pairs for english, 2023. URL
980 <https://arxiv.org/abs/1912.00582>.
- 981
- 982 Zhengzhuo Xu, Bowen Qu, Yiyang Qi, Sinan Du, Chengjin Xu, Chun Yuan, and Jian Guo. Chart-
983 moe: Mixture of diversely aligned expert connector for chart understanding. *arXiv preprint*
984 *arXiv:2409.03277*, 2024.
- 985 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu
986 Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin,
987 Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen.
988 Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for
989 expert agi. *ArXiv*, abs/2311.16502, 2023. URL [https://api.semanticscholar.org/
990 CorpusID:265466525](https://api.semanticscholar.org/CorpusID:265466525).
- 991 S E Yuksel, J N Wilson, and P D Gader. Twenty Years of Mixture of Experts. *IEEE Transactions on*
992 *Neural Networks and Learning Systems*, 23(8):1177–1193, 2012. ISSN 2162-2388 VO - 23. doi:
993 10.1109/TNNLS.2012.2200299.
- 994
- 995 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine
996 really finish your sentence?, 2019. URL <https://arxiv.org/abs/1905.07830>.
- 997 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language
998 image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
999 pp. 11975–11986, 2023.
- 1000
- 1001 Mian Zhang, Xianjun Yang, Xinlu Zhang, Travis Labrum, Jamie C. Chiu, Shaun M. Eack, Fei Fang,
1002 William Yang Wang, and Zhiyu Zoey Chen. Cbt-bench: Evaluating large language models on
1003 assisting cognitive behavior therapy, 2025a. URL <https://arxiv.org/abs/2410.13218>.
- 1004 Yi-Fan Zhang, Huanyu Zhang, Haochen Tian, Chaoyou Fu, Shuangqing Zhang, Junfei Wu, Feng Li,
1005 Kun Wang, Qingsong Wen, Zhang Zhang, Liang Wang, Rong Jin, and Tieniu Tan. Mme-realworld:
1006 Could your multimodal llm challenge high-resolution real-world scenarios that are difficult for
1007 humans?, 2025b. URL <https://arxiv.org/abs/2408.13257>.
- 1008 Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai,
1009 zhifeng Chen, Quoc V Le, and James Laudon. Mixture-of-Experts with Expert Choice Rout-
1010 ing. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Ad-
1011 vances in Neural Information Processing Systems*, volume 35, pp. 7103–7114. Curran Asso-
1012 ciates, Inc., 2022. URL [https://proceedings.neurips.cc/paper_files/paper/
1013 2022/file/2f00ecd787b432c1d36f3de9800728eb-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/2f00ecd787b432c1d36f3de9800728eb-Paper-Conference.pdf).
- 1014
- 1015 Simiao Zuo, Xiaodong Liu, Jian Jiao, Young Jin Kim, Hany Hassan, Ruofei Zhang, Jianfeng Gao,
1016 and Tuo Zhao. Taming Sparsely Activated Transformer with Stochastic Experts. In *International
1017 Conference on Learning Representations*, 2022. URL [https://openreview.net/forum?
1018 id=B72HXs80q4](https://openreview.net/forum?id=B72HXs80q4).
- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

Supplement to “CompeteSMoE – Statistically Guaranteed Mixture of Experts Training via Competition”

This document provides the supplementary materials for the paper CompeteSMoE – Statistically Guaranteed Mixture of Experts Training via Competition, and is organized as follows.

Contents

A Summary of Main Notations	21
B Adaptive Layer-wise Competition Control	22
C Ablation Study	22
C.1 Effect of Component-wise Design on Model Performance.	22
C.2 Joint Optimization of Task Loss and Competition Policy	23
C.3 Effectiveness of Activation Functions in the Competition Mechanism	24
C.4 Evaluation of Mean and Norm Strategies for Competition Mechanism	24
C.5 Evaluation of Distillation Loss Effectiveness	25
D Evolutionary Analysis of CompeteSMoE Behavior	25
D.1 Activation Magnitude Dynamics Across Layers	25
D.2 Evolution of Routing Score Dynamics During Training	26
D.3 Evolution Analysis of CompeteSMoE Across Multiple Sparse-Expert Configurations	26
E Further Analysis of Router Behavior	27
F Additional Experimental Results	28
G Additional A_{\max} Sensitivity Results	28
H Experimental Setup Details	29
H.1 Vision-Language Model (VLM)	29
Training Stages.	29
Architecture.	29
Hyperparameters.	29
Dataset.	30
H.2 Language Model Pretraining	30
Architecture.	30
Hyperparameters.	30
Dataset.	30
H.3 CompeteSMoE Configuration	31
Hyperparameters.	31
I Training Curves on Vision-Language Benchmarks	31
J Additional Theoretical Results	31
K Proof of Theoretical Results	32
K.1 Proof of Theorem 3.2	32
K.2 Proof of Theorem J.1	39
K.3 Proof of Proposition 3.1	45
K.3.1 Preliminaries	45
K.3.2 Main Proof	47
L LLM Usage	48

M Broader Impact

48

A SUMMARY OF MAIN NOTATIONS

We summarize the main notations used in the main paper in Table 9, including those introduced later in the supplementary material.

Table 9: Summary of Main Notations.

Symbol	Description
\mathcal{R}, W_r	Router network (function) and its parameter
g, W_e	Expert network (function), and its parameter
\mathbf{x}	Input
$\mathbf{s}, \mathbf{s}_R, \mathbf{s}_C$	Affinity scores, affinity scores from the router, affinity scores from competition
$\text{TopK}_{-\infty}$	Function retaining the K largest vector elements and setting others to $-\infty$
TopK_0	Function retaining the K largest vector elements and setting others to 0
K	Number of experts activated per input
N	The total number of experts
$[M]$	Set of $\{1, 2, \dots, M\}$ for any positive integer M
\hat{y}, y	Predicted output, ground truth
t	Current t -th iteration
T	Total number of training steps
l	The l -th SMOE layer
L	Total number of SMOE layers in the model
κ	Activation function
σ	Scoring function
$\mathbb{E}[\cdot]$	Mean of vector elements
e	Base of the exponential function
I_C	Indices of experts who won in the competition mechanism
α	Hyper-parameter prioritizing winning experts in distillation loss
γ	Hyper-parameter for distillation loss
β	Hyper-parameter for diversity loss
ω	Bernoulli probability for scheduling competition in each layer
A_{\max}	Maximum number of layers that can perform competition on a single time step
$\lambda(t)$	A scheduler determining whether to perform competition at the t -th step
$\Lambda(l)$	A vector storing the results of the scheduler $\lambda(t)$ at all time steps of the l -th layer
\mathcal{L}_{NLL}	Negative log-likelihood function (task loss)
$\mathcal{L}_{\mathcal{D}}$	Distillation loss
\mathcal{L}_{div}	Diversity loss
ξ_t	Step size
\mathcal{D}	A benchmark dataset for evaluation
Q_{prev}	Cumulative competition activations over layers 1 to $l-1$
$a_n = \mathcal{O}(b_n)$ or $a_n \lesssim b_n$	If $a_n \leq Cb_n$ for all $n \in \mathbb{N}$, where $C > 0$ is some universal constant
$a_n = \mathcal{O}_P(b_n)$	$\forall \epsilon > 0, \exists M > 0 : \mathbb{P}(A_n/b_n > M) < \epsilon$ for all sufficiently large n
$a_n = \tilde{\mathcal{O}}_P(b_n)$	$a_n = \mathcal{O}_P(b_n \log^c(b_n))$, for some $c > 0$.
$w^{(u)}, w_u$	The u -th entry of a vector $w \in \mathbb{R}^d$
w^z	$w^z = w_1^{z_1} w_2^{z_2} \dots w_d^{z_d}$, for any vector $w \in \mathbb{R}^d$ and $z \in \mathbb{N}^d$
$ w $	$ w := w_1 + w_2 + \dots + w_d$, for any vector $w \in \mathbb{R}^d$
$z!$	$z! := z_1! z_2! \dots z_d!$, for any vector $z \in \mathbb{N}^d$
N^*	The number of ground-truth experts
$f(\cdot \mu, \nu)$	Univariate Gaussian density with mean μ and variance ν
G_*	Ground-truth mixing measure
δ	Dirac measure
m	Lebesgue measure
Θ	Parameter space
d_1	Dimension of input space
d_2	Dimension of expert parameter space
\hat{G}_n	Maximum likelihood estimator for G_*
$\ \cdot\ , \ \cdot\ _1$	ℓ_2 -norm and ℓ_1 -norm value
$ A $	Cardinality of any set A
$h(p_1, p_2)$	Hellinger distance $h(p_1, p_2) := \left(\frac{1}{2} \int (\sqrt{p_1} - \sqrt{p_2})^2 dm \right)^{1/2}$ for any densities p_1, p_2
$V(p_1, p_2)$	Total Variation distance $V(p_1, p_2) := \frac{1}{2} \int p_1 - p_2 dm$ for any densities p_1, p_2

B ADAPTIVE LAYER-WISE COMPETITION CONTROL

While scheduled training reduces computational overhead, excessive simultaneous competition activations across multiple SMoE layers can destabilize the training process. To address this, we propose a dynamic mechanism that regulates the number of active competition layers at each training step, enhancing training efficiency. This is achieved by enforcing a global constraint on the maximum number of simultaneously active layers.

For a given layer l , we compute the cumulative competition activations from all preceding layers (i.e., layers 1 through $l - 1$) as:

$$Q_{\text{prev}} = \sum_{i=1}^{l-1} \Lambda(i), \quad (8)$$

where $\Lambda(i) \in \mathbb{R}^T$ denotes the activation state vector of layer i over T training steps, and $Q_{\text{prev}} \in \mathbb{R}^T$ represents the cumulative competition activations up to layer $l - 1$.

A predefined threshold $A_{\text{max}} \in \mathbb{R}$ governs the total number of active layers permitted per training step. If activating layer l at step t exceeds this threshold i.e., if $Q_{\text{prev}}(t) + \Lambda(l, t) > A_{\text{max}}$ with $\Lambda(l, t) = 1$ we redistribute the activation to an alternative step $t' \neq t$ satisfying:

$$Q_{\text{prev}}(t') + 1 \leq A_{\text{max}}, \quad t' \in \{1, \dots, T\}, \quad \Lambda(l, t') = 0. \quad (9)$$

Upon identifying t' , we update the activation schedule by setting $\Lambda(l, t') = 1$ and $\Lambda(l, t) = 0$. Empirical results indicate that only 0% to 7% of layers are active at any step, ensuring the availability of suitable t' satisfying Eq. 9.

In summary, this approach dynamically balances competition activations across layers, substantially reducing computational overhead while maintaining training stability for CompeteSMoE. Notably, the value of A_{max} depends on several factors such as model architecture, batch size, and available GPU memory, and may vary if the experiments are conducted in a different environment.

C ABLATION STUDY

We conducted an ablation study on a 5.1B parameter VLM, evaluating performance across various configurations. The best performance was observed with the large-scale model.

C.1 EFFECT OF COMPONENT-WISE DESIGN ON MODEL PERFORMANCE.

Table 10: Comprehensive component ablation study of CompeteSMoE across nine benchmarks.

Method	Scheduler	Competition Mechanism	Diversity Loss	Dense MoE	AVG Acc \uparrow	AVG Rank \downarrow
CompeteSMoE	✓	✓	✓	✗	53.21	1.78
(1)	✓	✓	✗	✗	52.90	3.11
(2)	✓	✗	✓	✗	52.71	3.89
(3)	✓	✗	✗	✓	52.70	4.11
(4)	✓	✗	✓	✓	52.79	3.44
SMoE	✗	✗	✗	✗	52.47	4.67

As shown in Table 10, we conduct a component-wise ablation study of the proposed CompeteSMoE model across nine benchmark datasets. Both the Competition Mechanism (1) and Diversity Loss (2) independently yield improvements over the standard SMoE baseline. Specifically, disabling DL results in a 0.49% drop in average accuracy, while removing CM leads to a smaller degradation of 0.30%. These findings suggest that CM contributes more significantly to overall model performance than DL when assessed in isolation.

In addition, we introduce two extended variants, models (3) and (4), inspired by the DenseMoE design (Pan et al., 2024), which activate all experts to compute the output but only occasionally during training using the same scheduler configuration. Interestingly, CompeteSMoE still consistently

outperforms. While dense activation provides a modest improvement over vanilla SMOE, it remains inferior to CompeteSMoE. This indicates that the performance gain arises not from dense expert activation per se, but from the competitive dynamics introduced by CM.

Crucially, in DenseMoE like variants (3) and (4), all experts are activated during both the forward and backward passes, leading to significantly increased computational cost. In contrast, CompeteSMoE activates all experts only to compute affinity scores, and then selects only the top-K winning experts to contribute to the final output and receive gradients. This hybrid mechanism enables more effective routing supervision while maintaining the computational efficiency characteristic of sparse MoE models.

C.2 JOINT OPTIMIZATION OF TASK LOSS AND COMPETITION POLICY

Table 11: Ablation study showing the impact of task optimization and competition policy matching on performance across 9 benchmarks.

Model	Task Loss	Match Competition Policy	AVG Acc	AVG Rank
CompeteSMoE	✓	✓	53.21	1.33
CompeteSMoE – Competition Policy Only	✗	✓	52.84	2.11
SMoE	✓	✗	52.47	2.56

Our analysis in Section 3 established that the competition policy alone is theoretically sufficient to achieve faster convergence compared to vanilla SMOE. However, in practice, CompeteSMoE stacks multiple SMOE layers, each independently equipped with a competition mechanism. Such a deep architecture requires significantly more training samples for convergence, often exceeding the dataset sizes available and making training infeasible under our hardware constraints. Therefore, we jointly optimize for the task loss and match the competition policy to facilitate practical and efficient learning.

The two supervision signals play complementary roles. When competition is inactive, the task loss gradient updates the router by adjusting affinity scores for the selected experts only, since inactive experts receive no gradients. When competition is active, its gradient provides updates for all experts, including those not selected in the final prediction. Combining both objectives thus provides more robust supervision and accelerates learning.

To validate this intuition, we conducted an ablation study on 9 benchmarks, isolating the effect of each supervision signal. Results are shown in Table 11. Jointly optimizing both signals yields the best average accuracy and rank. Interestingly, training the router solely to match the competition policy without any task loss supervision already surpasses the standard SMOE. This demonstrates that competition driven learning alone is capable of discovering stronger routing policies, even though the competition policy is active in only 7% of training steps. Despite such sparse updates, the router still learns significantly better expert selection than SMOE, indicating that the competition policy acts as a strong inductive bias.

Finally, the key properties of competition remain central in CompeteSMoE. As shown in Figure 1, CompeteSMoE achieves both faster and stronger convergence: after 8 hours of training, it already outperforms baselines trained for up to 14 hours. Moreover, Section 5.2.2 a) demonstrates that when deviating from the learned router policy (replacing the Top-1 expert with the Top-($K+1$) expert), SMOE surprisingly improves, whereas CompeteSMoE degrades. This indicates that the routing policy learned by SMOE is suboptimal, while CompeteSMoE produces a stronger and more consistent routing strategy.

C.3 EFFECTIVENESS OF ACTIVATION FUNCTIONS IN THE COMPETITION MECHANISM

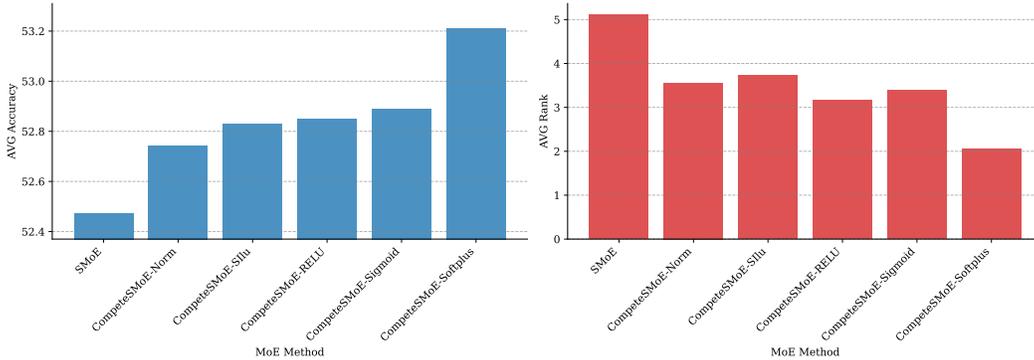


Figure 4: Performance comparison of different activation functions used within the Competition Mechanism across 9 benchmarks.

This section investigates how different activation functions influence the effectiveness of the Competition Mechanism. Specifically, we examine their role in computing expert affinity scores, originally defined in Eq. 2.2. To support a broader class of diversity-inducing functions, we generalize the affinity score formulation as follows:

$$s_i = \mathbb{E}[\kappa(g(\mathbf{x}, W_{e_i}))], \quad \forall i \in [N], \quad (10)$$

where $\kappa(\cdot)$ denotes the activation function applied to the expert response. This generalization allows the Competition Mechanism to flexibly incorporate a variety of activation profiles for expert selection. As shown in Figure 4, we compare several widely used activation functions within this framework, including Softplus, SiLU, Sigmoid, ReLU, and Softmax. Among them, Softplus consistently yields the highest average accuracy and ranking across tasks. We attribute this to its smooth and well behaved response curve, which softly suppresses negative values while preserving the magnitude of positive inputs. This behavior enables it to retain informative signals across the activation range, maintaining both representational richness and continuous gradient flow two properties critical for stable optimization. By contrast, Sigmoid compresses the entire input domain into the $[0, 1]$ interval, which can lead to vanishing gradients and loss of signal, especially for inputs with large magnitude. ReLU, although preserving positive values, entirely discards negative activations, potentially eliminating useful information. SiLU and Softmax lie between these extremes but still fall short of the balance offered by Softplus. We also explored an alternative formulation using the exponential function: $\mathbb{E}[e^{g(\mathbf{x}, W_{e_i})}]$. However, this variant led to uncontrolled growth in output magnitudes, causing numerical instability and NaN values during training. In contrast, Softplus provides a smooth approximation to the exponential function while mitigating such instability, making it a more robust choice for this setting.

In summary, activation functions that gently suppress negative activations while maintaining near-linear behavior on the positive side such as Softplus are better aligned with the needs of the Competition Mechanism. Their balanced characteristics lead to more stable expert affinity computation and improved end-task performance.

C.4 EVALUATION OF MEAN AND NORM STRATEGIES FOR COMPETITION MECHANISM

We conduct an empirical investigation to compare the mean-based strategy, as defined in Eq. 2.2, with a norm-based formulation. Specifically, we compute the affinity score of expert i using the L2 norm of its output vector:

$$s_i = \|g(\mathbf{x}, W_{e_i})\|, \quad \forall i \in [N], \quad (11)$$

As shown in Figure 4, the CompeteSMoE-Norm variant using Equation 11 yields higher performance compared to the SMoE standard. However, when we switch to the CompeteSMoE-Softplus configuration that employs a mean based strategy, a substantial improvement is observed in both average accuracy and ranking. In conclusion, the mean-based strategy proves to be the most effective setting for expert output aggregation within the Competition Mechanism.

C.5 EVALUATION OF DISTILLATION LOSS EFFECTIVENESS

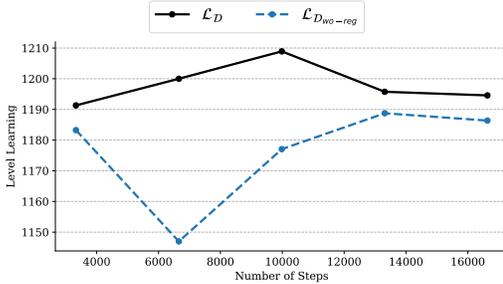


Table 12: Performance comparison between \mathcal{L}_D and $\mathcal{L}_{D_{wo-reg}}$ across 9 benchmark datasets.

Loss Function	Avg. Acc	Avg. Rank
$\mathcal{L}_{D_{wo-reg}}$	52.92	1.78
\mathcal{L}_D	53.21	1.22

Figure 5: Learning performance of \mathcal{L}_D and $\mathcal{L}_{D_{wo-reg}}$ measured by the Level Learning metric at every 20% of training steps on the MMBench-EN benchmark.

In Section 3, we established the theoretical foundation for the competition mechanism and demonstrated its empirical effectiveness in Table 1. A key challenge in optimizing the router network is accurately modeling the distribution of competitive routing decisions. We carefully investigated two objective functions: the distillation loss \mathcal{L}_D (see details in Eq. 1) and a variant distillation loss $\mathcal{L}_{D_{wo-reg}}$ without the regularization term, which emphasizes penalizing experts who won the competition. We define $\mathcal{L}_{D_{wo-reg}}$ as follows:

$$\mathcal{L}_{D_{wo-reg}}(s_R, s_C) = \text{MSE}(s_R, s_C) \tag{12}$$

Figure 5 illustrates the progression of the Level Learning (LL) metric, which measures the average number of Top- K experts selected by the router network that align with the Top- K experts from the competition mechanism. A high LL value indicates that the router network effectively learns from the competition mechanism, whereas a low value suggests poor learning performance. Notably, \mathcal{L}_D consistently enables faster and more stable convergence compared to $\mathcal{L}_{D_{wo-reg}}$. In particular, during the initial 60% of training (up to 9,600 steps), \mathcal{L}_D maintains a clear advantage, effectively mitigating the early performance drop observed with $\mathcal{L}_{D_{wo-reg}}$. Moreover, \mathcal{L}_D achieves a peak LL score of 1210 by 12,000 steps, surpassing the $\mathcal{L}_{D_{wo-reg}}$ peak of 1190, and exhibits more stable learning dynamics in later stages.

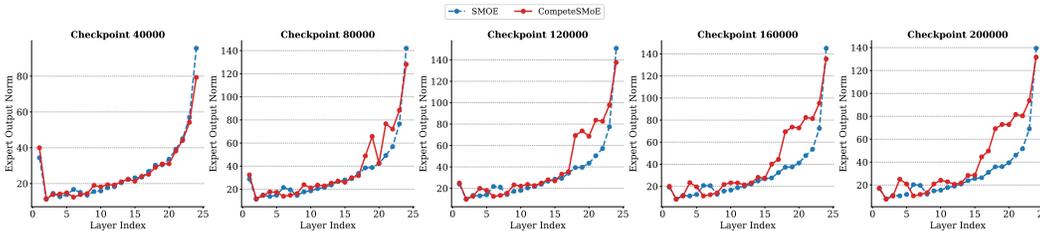
Additionally, quantitative results in Table 12 further confirm this trend, with \mathcal{L}_D yielding a higher average accuracy (53.21% vs. 52.92%) and a lower average rank (1.22 vs. 1.78) across nine benchmarks. These findings underscore the effectiveness of \mathcal{L}_D in guiding the router network to better approximate the competition mechanism. Furthermore, they suggest its potential as a preferred optimization objective in competitive MoE architectures.

D EVOLUTIONARY ANALYSIS OF COMPETESMOE BEHAVIOR

D.1 ACTIVATION MAGNITUDE DYNAMICS ACROSS LAYERS

In Figure 6, we analyze how expert output norms evolve across layers and training checkpoints for both SMoE and CompeteSMoE, following the norm-based methodology in MoEUT (Csordás et al., 2024). Although CompeteSMoE exhibits slightly larger activations, this increase is highly localized to the middle layers (approximately Layers 16–25). The early layers and the final layers

1350
1351
1352
1353
1354
1355
1356
1357



1358
1359
1360
1361

Figure 6: Layer-wise expert output norms across training progress on a 1B language pretraining model.

1362
1363
1364
1365
1366

remain largely unchanged, and several of the deepest layers even show smaller activation magnitudes than SMOE. These trends indicate that the competition mechanism does not induce a monotonic or runaway growth pattern. Overall, the activation changes introduced by CompeteSMoE are moderate, layer-specific, and non-accumulative, suggesting that the method maintains stable activation scales across the network and avoids globally inflated magnitudes.

1367

1368

D.2 EVOLUTION OF ROUTING SCORE DYNAMICS DURING TRAINING

1369

1370

1371

1372

1373

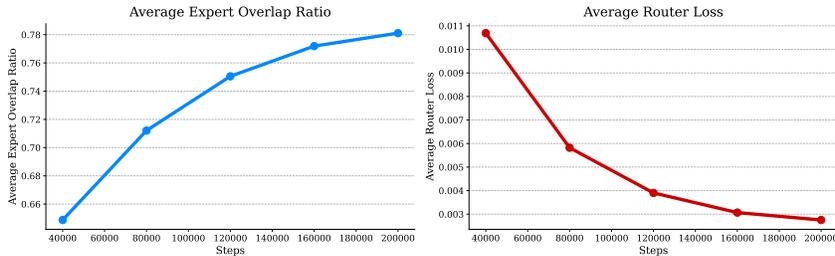
1374

1375

1376

1377

1378



1379
1380
1381
1382

Figure 7: Evolution of router learning dynamics during training in the CompeteSMoE algorithm. The left plot presents the **expert overlap ratio** between the router’s Top- K selections and the competition mechanism, while the right plot shows the **router distillation loss**, which quantifies the discrepancy between router scores and competition scores.

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

D.3 EVOLUTION ANALYSIS OF COMPETESMOE ACROSS MULTIPLE SPARSE-EXPERT CONFIGURATIONS

1396

1397

1398

1399

1400

1401

1402

1403

We evaluate the robustness of CompeteSMoE under varying levels of expert sparsity. As illustrated in Figure 8, we train a 0.3B-parameter language model pretrain on a 13B-token corpus and perform a three-level sparsity ablation with $K \in \{2, 4, 8\}$ under an $N = 24$ -expert architecture. This setting enables us to isolate the effect of expert activation density on model performance. Across all sparsity configurations, CompeteSMoE consistently outperforms the standard SMOE baseline throughout the training trajectory. Importantly, the improvements are stable and do not depend on any particular selection of K . This pattern indicates that the performance gains arise from the competition-based routing mechanism itself, rather than from architectural idiosyncrasies or favorable hyperparameter choices.

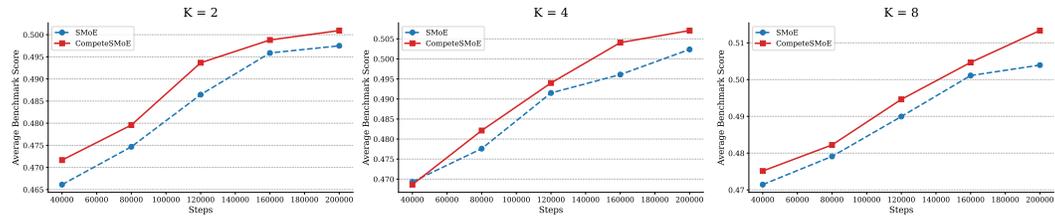


Figure 8: Comparison of average performance across 8 benchmarks between SMoE and CompeteSMoE under expert configurations $K \in \{2, 4, 8\}$ and an $N = 24$ -expert architecture, evaluated over training steps.

Overall, these results demonstrate that CompeteSMoE provides a reliable performance advantage across sparse-expert regimes. The method remains effective even when the degree of sparsity varies substantially, underscoring that its benefits stem from genuine algorithmic advances rather than from artifacts of model configuration, dataset size, or training dynamics.

E FURTHER ANALYSIS OF ROUTER BEHAVIOR

In this section, we further analyze about router behavior in SMoE and CompeteSMoE.

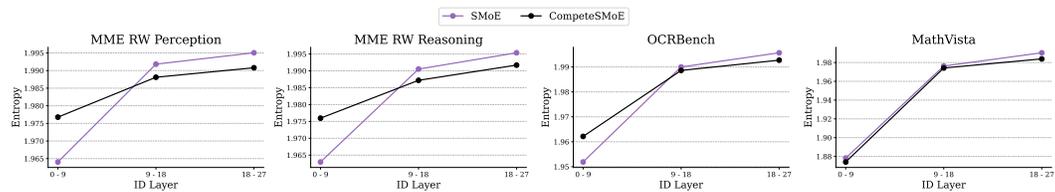


Figure 9: Entropy analysis of expert selection frequency across perception and reasoning tasks. Lower entropy indicates higher specialization in expert routing.

(a) Experts distribution on Reasoning and Perception. As illustrated in Figure 9, we analyze the entropy of expert distribution across layers for SMoE and CompeteSMoE algorithms, evaluated on three benchmarks: MME Real-World Perception and OCR Bench for perception capacity, and MME Real-World Reasoning and MathVista for reasoning capacity. On perception tasks, CompeteSMoE exhibits higher entropy in the early layers, indicating exploratory behavior, but significantly reduces entropy in the middle and final layers. In contrast, on MathVista a benchmark requiring higher-level reasoning CompeteSMoE maintains low entropy in the early and intermediate layers, approaching entropy levels similar to SMoE in the final layers. Both models demonstrate increasing entropy toward the final layers, suggesting more balanced expert allocation as the network deepens, consistent with typical Transformer-based architectures where later layers aggregate information from multiple upstream experts. Regarding the **representation collapse issue**, both SMoE and CompeteSMoE achieve a high degree of balance in expert distribution, with entropy scores exceeding 1.99 (compared to the maximum entropy of 2 for four experts).

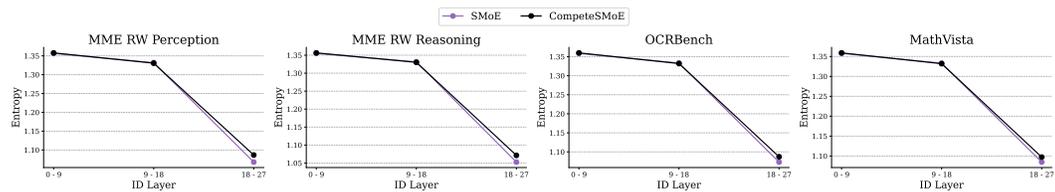


Figure 10: Layer-wise entropy of expert weight distributions for CompeteSMoE and SMoE across three tasks: Real-World Perception, Real-World Reasoning, and Mathematical Reasoning.

(b) **Effective Expert Aggregation via Weight Distribution.** As shown in Figure 10, we analyze the entropy of expert weight distributions across layers and tasks, which reflects how expert contributions are aggregated. Lower entropy typically suggests more confident expert selection. Both SMoE and CompeteSMoE exhibit decreasing entropy across layers, implying increased decisiveness in expert routing at deeper layers. While SMoE generally maintains lower entropy, especially on MathVista, it tends to concentrate weights heavily on a small subset of experts. In contrast, CompeteSMoE distributes weights more evenly among the selected experts. This balanced aggregation allows CompeteSMoE to better leverage complementary knowledge from multiple experts. Finally, we observe a slight difference between the two models, with both showing a trend toward more confident weight distributions in the final layers.

F ADDITIONAL EXPERIMENTAL RESULTS

Table 13: Performance comparison of SMoE strategies in the VIT pretraining setting, where only the MLP connectors are initialized as SMoE layers from scratch, using a $\approx 4\text{B}$ -parameter model. **Bold** values indicate the best result, while underlined values denote the second best. \uparrow / \downarrow indicate higher/lower is better.

Method	AI2D	Text VQA	GQA	MM Bench	Hallusion	Math Vista	MMMU	MMStar	POPE	OCR	MME RWL	Avg. Acc \uparrow	Avg. Rank \downarrow
SMoE	59.84	39.56	56.29	66.32	46.69	29.10	38.11	37.55	86.14	32.40	32.86	47.72	4.36
XMoE	<u>61.59</u>	39.40	55.92	65.98	45.32	<u>29.80</u>	39.00	<u>39.82</u>	86.58	32.10	31.21	47.88	3.86
PCosine	56.09	39.61	49.98	52.41	40.59	27.00	38.78	37.38	86.20	22.70	30.22	43.72	6.27
σ -MoE	<u>61.59</u>	39.18	56.60	65.98	44.80	30.00	39.00	37.89	86.28	32.10	31.21	47.69	4.14
SharedE-V2	61.06	39.20	56.14	64.00	46.58	29.30	39.11	39.06	87.13	31.80	33.25	47.88	3.91
SharedE-V3	61.37	39.60	<u>56.33</u>	<u>66.58</u>	45.01	29.30	<u>37.56</u>	39.12	<u>86.97</u>	<u>32.20</u>	33.13	<u>47.92</u>	<u>3.41</u>
CompeteSMoE	61.76	39.71	56.00	67.61	46.79	29.30	39.56	40.12	86.36	<u>32.20</u>	34.18	48.51	2.05

In Table 13, we report additional results using a vision-language model where only the MLP connectors are replaced with SMoE layers consisting of 8 experts, with 4 experts active per token. The vision encoder and language model (LLM) are kept dense and frozen during training, while the MLP connectors are unfrozen, following the setup described in Xu et al. (2024). These SMoE layers are trained from scratch using the same dataset and training configuration as the VIT stage, enabling a controlled analysis of sparse upcycling in isolation. Under this setup, CompeteSMoE consistently outperforms all baseline methods.

G ADDITIONAL A_{\max} SENSITIVITY RESULTS

To systematically assess the sensitivity of A_{\max} across a broader range of values, we conducted ablation studies by pretraining a smaller 0.3B-parameter language model using $A_{\max} \in \{3, 6, 9\}$, with $N = 24$ and $K = 8$ experts, trained on 13B tokens. The results are summarized in Table 14.

Table 14: Comparison of CompeteSMoE performance across maximum activation thresholds ($A_{\max} \in \{3, 6, 9\}$).

Method	SMoE	CompeteSMoE		
		$A_{\max} = 3$	$A_{\max} = 6$	$A_{\max} = 9$
Avg. across 8 benchmarks	50.39	50.95	51.33	51.11

As shown in Table 14, the performance of CompeteSMoE steadily increases as A_{\max} rises, peaking at $A_{\max} = 6$ with a clear improvement from 3 to 6, then slightly decreases as A_{\max} increases from 6 to 9, although this fluctuation is negligible. This trend indicates that larger A_{\max} values generally result in better and more stable performance, and the differences between configurations gradually narrow, reflecting greater robustness. Notably, all CompeteSMoE settings with $A_{\max} \in \{3, 6, 9\}$ outperform the SMoE baseline. Similarly, we also adopted $A_{\max} = 6$ for the 1B language model pretraining experiment in Table 2, where it again achieved the best results. These findings further reinforce our recommendation regarding the choice of A_{\max} .

1512 H EXPERIMENTAL SETUP DETAILS

1513 H.1 VISION-LANGUAGE MODEL (VLM)

1514 **Training Stages.** We adopt a three-stage training pipeline inspired by prior works (Nguyen et al.,
1517 2024c; Li et al., 2024), designed to incrementally adapt and integrate the vision and language
1518 modalities for multimodal instruction tuning. Table 15 summarizes the training status of each model
1519 component throughout the stages.

1520 Table 15: Component Training States at Each Stage

1522 Stage	1523 LLM	1523 MLP Connector	1523 Vision Encoder
1524 Pre-Training	1524 Frozen	1524 Trainable	1524 Frozen
1525 Pre-Finetuning	1525 Trainable	1525 Trainable	1525 Trainable
1526 Visual Instruction Tuning	1526 Trainable	1526 Trainable	1526 Trainable

- 1528 • **(1) Pre-Training (PT):** In the first stage, only the MLP connector is trained, while the vision
1529 encoder and the language model (LLM) are kept frozen. This stage focuses on aligning visual
1530 features with the language embedding space, establishing a stable initialization without
1531 perturbing the frozen backbones.
- 1532 • **(2) Pre-Finetuning (PFT):** All components including the vision encoder, MLP connector,
1533 and LLM are unfrozen and trained jointly using a dense architecture. This warm-up stage
1534 strengthens the cross-modal representation and stabilizes the model before introducing
1535 sparsity.
- 1536 • **(3) Visual Instruction Tuning (VIT):** In the final stage, we apply Sparse Upcycling (Ko-
1537 matsuzaki et al., 2023) by replacing selected MLP layers in both the vision encoder and the
1538 connector with Top- K sparsely gated MoE blocks. As this setup has become the standard
1539 practice in recent vision-language research (Nguyen et al., 2024c; Li et al., 2024; Shu et al.,
1540 2024; Lin et al., 2024), evaluating MoE algorithms under this setting offers more meaningful
1541 and practically relevant comparisons. Each expert is initialized from its corresponding
1542 pretrained MLP, while the Top- K router is learned from scratch. To promote balanced
1543 expert utilization, we apply standard auxiliary objectives, including load-balancing loss and
1544 z-loss. All components remain fully trainable during this stage, and all compared methods
1545 are trained and evaluated under this unified VIT setup.

1547 **Architecture.** We adopt a modular design for the VIT-stage model, composed of a vision encoder,
1548 an MLP-based connector, and a pretrained language model backbone. The detailed architecture and
1549 MoE configuration are summarized in Table 16. During this stage, MoE layers are applied to the
1550 vision encoder and connector, while the language model remains dense.

1551 Table 16: Architecture and parameter breakdown for each component in the 5.1B VIT-stage model,
1552 with MoE usage indicated.

1554 Component	1554 Version / Variant	1554 Parameters	1554 SMoE
1556 Vision Encoder	1556 SigLIP-SO400M-Patch14-224	1.20B	✓
1557 MLP Connector	1557 -	66M	✓
1558 Language Model	1558 Phi-3.5 Mini Instruct	3.82B	✗
1559 Total	1559 -	1559 5.1B	1559 -

1562 **Hyperparameters.** Table 17 lists the training hyperparameters for each stage. During VIT, SMoE
1563 layers are introduced by upcycling dense MLPs in the vision encoder and MLP connector. Each
1564 SMoE block contains $N_E = 4$ experts, with $K = 2$ experts selected per token. Given the sensitivity
1565 of large-scale model training to initialization and randomness, we ensure that all baseline models
share identical starting conditions. Specifically, the router networks within MoE blocks are initialized

from scratch using Gaussian noise $\mathcal{N}(0, 0.02)$, following the scheme used in the official GPT-2 implementation.¹ To guarantee reproducibility, we fix the random seed to 42 across all experiments.

Table 17: Hyperparameter configurations for the three training stages of Phi-3.5 Mini.

Hyperparameter	PT	PFT	VIT
Learning rate	1e-3	2e-6	4e-6
Schedule	Cosine	Cosine	Cosine
Batch size / GPU	64	6	5
GPUs	4×H100	4×H100	4×H100
ZeRO stage	ZeRO-2	ZeRO-2	ZeRO-3
Optimizer	AdamW	AdamW	AdamW
MoE blocks	No	No	Yes
Balance loss coeff.	0.0	0.0	0.01
Z-loss coeff.	0.0	0.0	0.001
Max sequence length	2048	2048	2048

Dataset. The full pipeline uses over 1B tokens, spanning: (1) LCS-558K (Liu et al., 2023a) (pretraining), (2) ALLaVA-Caption (Chen et al., 2024a) 708K (pre-finetuning), and (3) LLaVA-665K (instruction tuning) (Liu et al., 2024a), consistent with LibMoE (Nguyen et al., 2024c) and CuMo (Li et al., 2024).

H.2 LANGUAGE MODEL PRETRAINING

Architecture. We pretrain a 1B-parameter and language model following the *original Transformer architecture*, where SMoE layers are integrated exclusively within the MLP blocks. Each SMoE layer consists of $N_E=24$ experts, with $K=8$ experts activated per token. Our architectural design is inspired by prior configurations such as SmoLLM2-1.7B (Allal et al., 2025) and MoEUT (Csordás et al., 2024). The model specification is summarized in Table 18.

Table 18: Architecture configuration of the pretrained 1B language model.

#Params	n_{layers}	d_{model}	d_{expert}	H	d_{head}	N_E	K
1B	24	1024	512	32	128	24	8

Hyperparameters. We adopt the training configuration proposed in MoEUT Csordás et al. (2024), as detailed in Table 19. All model weights are initialized following the MoEUT initialization scheme, and a fixed random seed of 42 is set for reproducibility.

Table 19: Pretraining hyperparameters for MoE language model.

Learning Rate	Schedule	Batch size	GPUs	Optimizer	Balance Coeff.	N_{warmup}
0.00025	Cosine	64	4×H100	AdamW	0.01	4000

Dataset. We train our models on 13B tokens sampled from the SlimPajama corpus (Soboleva et al., 2023). Training is conducted for 200K steps under this data regime.

¹<https://github.com/openai/gpt-2>

H.3 COMPETESMOE CONFIGURATION

Hyperparameters. Table 20 provides the key hyperparameters for training CompeteSMoE on the 5.1B VLM model. We warm up the MoE layers for 5% of total steps before enabling competition. The parameter A_{\max} limits the number of concurrently active competition layers and is tuned for training stability. All runs use the same seed for fair comparison. Additional ablations on ω and α are included in Appendix C.

Table 20: CompeteSMoE hyperparameters on the 5.1B-parameter model.

Warm-up	ω	γ	α	β	A_{\max}
0.05	0.07	0.01	0.1	0.005	9

I TRAINING CURVES ON VISION-LANGUAGE BENCHMARKS

In Figure 11, we include additional training performance curves for 9 benchmarks, supplementing the results presented in Figure 1.

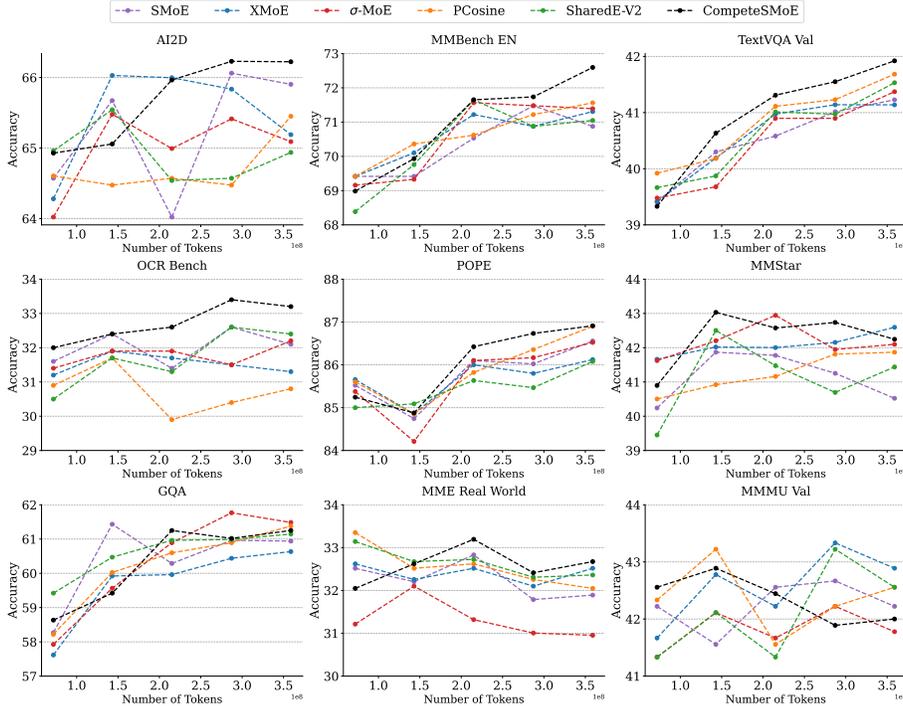


Figure 11: Training curves of CompeteSMoE compared to five advanced MoE algorithms on vision-language benchmarks.

J ADDITIONAL THEORETICAL RESULTS

In this appendix, we analyze the convergence behavior of Gaussian mixture of linear experts equipped with the competition mechanism. In particular, we consider experts of the linear form $g(X, (a, b)) := a^T X + b$, where $a \in \mathbb{R}^d$ and $b \in \mathbb{R}$. Then, the conditional density function $p_{G_*}(Y|X)$ in equation (3) becomes

$$p_{G_*}(Y|X) := \sum_{i=1}^{N^*} \frac{\exp(\log(1 + \exp((a_i^*)^T X + b_i^*)))}{\sum_{j=1}^{N^*} \exp(\log(1 + \exp((a_j^*)^T X + b_j^*)))} \cdot f(Y|(a_i^*)^T X + b_i^*, \nu_i^*). \quad (13)$$

Our ultimate goal is to compare the sample efficiency of this model to that without the competition mechanism (Nguyen et al., 2023a) in terms of expert estimation. For that purpose, we use a Voronoi loss tailored to the setting of linear experts, which is given by

$$\begin{aligned} \mathcal{L}_2(G, G_*) &:= \sum_{j=1}^{N^*} \left| \sum_{i \in \mathcal{C}_j} \exp(c_i) - \exp(c_j^*) \right| \\ &+ \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i) \left[\|a_i - a_j^*\| + |b_i - b_j^*| + |\nu_i - \nu_j^*| \right] \\ &+ \sum_{j \in [N^*]: |\mathcal{C}_j|>1} \sum_{i \in \mathcal{C}_j} \exp(c_i) \left[\|a_i - a_j^*\|^2 + |b_i - b_j^*|^2 + |\nu_i - \nu_j^*|^2 \right]. \end{aligned} \quad (14)$$

Equipped with the above Voronoi loss, we establish the convergence rate of parameter and expert estimations in the Gaussian mixture of linear experts with the competition in Theorem J.1.

Theorem J.1. *The following lower bound holds for any mixing measure $G \in \mathcal{G}_N(\Theta)$:*

$$\mathbb{E}_X[V(p_G(\cdot|X), p_{G_*}(\cdot|X))] \gtrsim \mathcal{L}_2(G, G_*). \quad (15)$$

This lower bound indicates that $\mathcal{L}_2(\hat{G}_n, G_*) = \mathcal{O}_P(\sqrt{\log(n)/n})$.

The proof of Theorem J.1 can be found in Appendix K.2. A few remarks regarding the results of this theorem are in order.

(i) *Parameter estimation rates.* The bound of the Voronoi loss $\mathcal{L}_2(\hat{G}_n, G_*)$ in Theorem J.1 reveals that the estimation rates for exact-specified parameters a_j^*, b_j^*, ν_j^* , i.e., for $j \in [N^*] : |\mathcal{C}_j| = 1$, are of parametric order $\tilde{\mathcal{O}}_P(n^{-1/2})$, whereas those for their over-specified counterparts, i.e., for $j \in [N^*] : |\mathcal{C}_j| > 1$, are slightly slower, of order $\tilde{\mathcal{O}}_P(n^{-1/4})$.

(ii) *Expert estimation rates.* Note that the input space is bounded, then we have

$$\left| (\hat{a}_i^n)^\top X + \hat{b}_i^n - (a_j^*)^\top X - b_j^* \right| \lesssim \|\hat{a}_i^n - a_j^*\| + |\hat{b}_i^n - b_j^*|,$$

for almost surely X . Consequently, the estimation rates for exact-specified and over-specified experts $(a_j^*)^\top X + b_j^*$ are also of orders $\tilde{\mathcal{O}}_P(n^{-1/2})$ and $\tilde{\mathcal{O}}_P(n^{-1/4})$, respectively.

(iii) *Sample efficiency of the competition mechanism.* Thus, we need polynomially many data points $\mathcal{O}(\epsilon^{-4})$ to estimate these linear experts with a given error $\epsilon > 0$. By contrast, when not using the competition mechanism Nguyen et al. (2023a), the linear expert estimation rates are substantially slowed down since they hinge on the solvability of some complex system of polynomial equations and are decelerated as the number of fitted experts grows. For example, if a linear expert $(a_j^*)^\top X + b_j^*$ is fitted by two experts (or three experts), that is, $|\mathcal{C}_j| = 2$ (or $|\mathcal{C}_j| = 3$), then the rate for estimating this linear expert is of order $\tilde{\mathcal{O}}_P(n^{-1/8})$ (or $\tilde{\mathcal{O}}_P(n^{-1/12})$). Therefore, we need $\mathcal{O}(\epsilon^{-8})$ (or $\mathcal{O}(\epsilon^{-12})$), to estimate this expert. For that reason, we claim that the Gaussian MoE becomes more sample-efficient when equipped with the competition mechanism.

K PROOF OF THEORETICAL RESULTS

K.1 PROOF OF THEOREM 3.2

In this proof, we aim to demonstrate that the following lower bound holds for any $G \in \mathcal{G}_N(\Theta)$:

$$\mathbb{E}_X[V(p_G(\cdot|X), p_{G_*}(\cdot|X))] \gtrsim \mathcal{L}_1(G, G_*). \quad (16)$$

For that purpose, we first establish the local part of the above bound, that is,

$$\lim_{\epsilon \rightarrow 0} \inf_{G \in \mathcal{G}_N(\Theta): \mathcal{L}_1(G, G_*) \leq \epsilon} \frac{\mathbb{E}_X[V(p_G(\cdot|X), p_{G_*}(\cdot|X))]}{\mathcal{L}_1(G, G_*)} > 0. \quad (17)$$

This local part implies that there exists a positive constant ϵ' that satisfies

$$\inf_{G \in \mathcal{G}_N(\Theta): \mathcal{L}_1(G, G_*) \leq \epsilon'} \frac{\mathbb{E}_X[V(p_G(\cdot|X), p_{G_*}(\cdot|X))]}{\mathcal{L}_1(G, G_*)} > 0.$$

Then, it is sufficient to derive the following global part of the bound in equation (16):

$$\inf_{G \in \mathcal{G}_N(\Theta): \mathcal{L}_1(G, G_*) > \varepsilon'} \frac{\mathbb{E}_X[V(p_G(\cdot|X), p_{G_*}(\cdot|X))]}{\mathcal{L}_1(G, G_*)} > 0. \quad (18)$$

Local part: In this part, we will establish the local part in equation (17) using the proof by contradiction method.

Suppose that the local part is not true, then we can find a sequence of mixing measures (G_n) given by $G_n := \sum_{i=1}^N \exp(c_i^n) \delta_{(W_{e_i}^n, \nu_i^n)} \in \mathcal{G}_N(\Theta)$ such that $\mathcal{L}_1(G_n, G_*) \rightarrow 0$ and

$$\mathbb{E}_X[V(p_{G_n}(\cdot|X), p_{G_*}(\cdot|X))]/\mathcal{L}_1(G_n, G_*) \rightarrow 0,$$

as $n \rightarrow \infty$. As we use asymptotic arguments in this proof, we may assume without loss of generality (WLOG) that the Voronoi cells $\mathcal{C}_j^n := \mathcal{C}_j(G_n)$ is independent of the sample size n . Then, the Voronoi loss of interest turns into

$$\begin{aligned} \mathcal{L}_1(G_n, G_*) &:= \sum_{j=1}^{N^*} \left| \sum_{i \in \mathcal{C}_j} \exp(c_i^n) - \exp(c_j^*) \right| + \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\|W_{e_i}^n - W_{e_j}^*\| + |\nu_i^n - \nu_j^*| \right] \\ &\quad + \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\|W_{e_i}^n - W_{e_j}^*\|^2 + |\nu_i^n - \nu_j^*|^2 \right]. \end{aligned} \quad (19)$$

Since $\mathcal{L}_1(G_n, G_*) \rightarrow 0$ as $n \rightarrow \infty$, we have $(W_{e_i}^n, \nu_i^n) \rightarrow (W_{e_j}^*, \nu_j^*)$ for all $j \in [N^*]$ and $i \in \mathcal{C}_j$.

Subsequently, we divide the rest of this proof into three main steps.

Step 1: Taylor expansion. In this step, we aim to decompose the term $T_n(Y|X) := \left[\sum_{j=1}^{N^*} \exp(\log(1 + \exp(g(x, W_{e_j}^*))) \right] \cdot [p_{G_n}(Y|X) - p_{G_*}(Y|X)]$ can be decomposed as

$$\begin{aligned} T_n(Y|X) &= \sum_{j=1}^{N^*} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp(g(X, W_{e_i}^n))) f(Y|g(X, W_{e_i}^n), \nu_i^n) \right. \\ &\quad \left. - \exp(\log(1 + \exp(g(X, W_{e_j}^*))) f(Y|g(X, W_{e_j}^*), \nu_j^*) \right] \\ &\quad - \sum_{j=1}^{N^*} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp(g(X, W_{e_i}^n))) - \exp(\log(1 + \exp(g(X, W_{e_j}^*))) \right] p_{G_n}(Y|X) \\ &\quad + \sum_{j=1}^{N^*} \left[\sum_{i \in \mathcal{C}_j} \exp(c_i^n) - \exp(c_j^*) \right] \cdot \exp(\log(1 + \exp(g(X, W_{e_j}^*))) [f(Y|g(X, W_{e_j}^*), \nu_j^*) - p_{G_n}(Y|X)] \\ &:= T_{n,1}(Y|X) - T_{n,2}(Y|X) + T_{n,3}(Y|X). \end{aligned}$$

Next, we continue to decompose the term $T_{n,1}(Y|X)$ as

$$\begin{aligned} T_{n,1}(Y|X) &= \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp(g(x, W_{e_i}^n))) f(Y|g(X, W_{e_i}^n), \nu_i^n) \right. \\ &\quad \left. - \exp(\log(1 + \exp(g(x, W_{e_j}^*))) f(Y|g(X, W_{e_j}^*), \nu_j^*) \right] \\ &\quad + \sum_{j \in [N^*]: |\mathcal{C}_j| > 1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp(g(x, W_{e_i}^n))) f(Y|g(X, W_{e_i}^n), \nu_i^n) \right. \\ &\quad \left. - \exp(\log(1 + \exp(g(x, W_{e_j}^*))) f(Y|g(X, W_{e_j}^*), \nu_j^*) \right] \\ &:= T_{n,1,1}(Y|X) + T_{n,1,2}(Y|X). \end{aligned}$$

Let us denote $F_\rho(Y|X; W_e, \nu) := \exp(\log(1 + \exp(g(X, W_e))) \frac{\partial^\rho f}{\partial g^\rho}(Y|g(X, W_e), \nu)$. By applying the first-order Taylor expansion to the function $F_0(Y|X; W_e, \nu)$ around the point $(W_{e_j}^*, \nu_j^*)$, we rewrite the term $T_{n,1,1}(Y|X)$ as

$$T_{n,1,1}(Y|X) = \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{\rho=0}^2 T_{n,1,1,\rho}^{(j)}(X) F_\rho(Y; X, W_{e_j}^*, \nu_j^*) + R_{n,1,1}(Y|X),$$

where $R_{n,1,1}(Y|X)$ is the Taylor remainder such that $R_{n,1,1}(Y|X)/\mathcal{L}_1(G_n, G_*) \rightarrow 0$ as $n \rightarrow \infty$, and

$$T_{n,1,1,0}^{(j)}(X) := \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \sum_{u=1}^{d_2} (\Delta W_{e_{ij}}^n)^{(u)} \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \cdot \frac{1}{1 + \exp(-g(X, W_{e_j}^*))},$$

$$T_{n,1,1,1}^{(j)}(X) := \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \sum_{u=1}^{d_2} (\Delta W_{e_{ij}}^n)^{(u)} \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*),$$

$$T_{n,1,1,2}^{(j)}(X) := \sum_{i \in \mathcal{C}_j} \frac{1}{2} \exp(c_i^n) (\Delta \nu_{ij}^n),$$

in which $\Delta W_{e_{ij}}^n := W_{e_i}^n - W_{e_j}^*$ and $\Delta \nu_{ij}^n := \nu_i^n - \nu_j^*$.

Meanwhile, by means of the second-order Taylor expansion, the term $T_{n,1,2}(Y|X)$ can be represented as

$$T_{n,1,2}(Y|X) = \sum_{j \in [N^*]: |\mathcal{C}_j| > 1} \sum_{\rho=0}^4 T_{n,1,2,\rho}^{(j)}(X) F_\rho(Y; X, W_{e_j}^*, \nu_j^*) + R_{n,1,2}(Y|X),$$

where $R_{n,1,2}(Y|X)$ is the Taylor remainder such that $R_{n,1,2}(Y|X)/\mathcal{L}_1(G_n, G_*) \rightarrow 0$ as $n \rightarrow \infty$, and

$$\begin{aligned} T_{n,1,2,0}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\sum_{u=1}^{d_2} (\Delta W_{e_{ij}}^n)^{(u)} \cdot \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} \right. \\ &\quad \left. + \sum_{u,v=1}^{d_2} \frac{(\Delta W_{e_{ij}}^n)^{(u)} (\Delta W_{e_{ij}}^n)^{(v)}}{1 + 1_{\{u=v\}}} \cdot \frac{\frac{\partial^2 g}{\partial W_e^{(u)} \partial W_e^{(v)}}(X, W_{e_j}^*) + \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} \right], \\ T_{n,1,2,1}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\sum_{u=1}^{d_2} (\Delta W_{e_{ij}}^n)^{(u)} \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \right. \\ &\quad \left. + \sum_{u,v=1}^{d_2} \frac{(\Delta W_{e_{ij}}^n)^{(u)} (\Delta W_{e_{ij}}^n)^{(v)}}{1 + 1_{\{u=v\}}} \left(\frac{2 \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} + \frac{\partial^2 g}{\partial W_e^{(u)} \partial W_e^{(v)}}(X, W_{e_j}^*) \right) \right], \\ T_{n,1,2,2}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\frac{1}{2} (\Delta \nu_{ij}^n) + \sum_{u,v=1}^{d_2} \frac{(\Delta W_{e_{ij}}^n)^{(u)} (\Delta W_{e_{ij}}^n)^{(v)}}{1 + 1_{\{u=v\}}} \cdot \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*) \right. \\ &\quad \left. + \sum_{u=1}^{d_2} (\Delta W_{e_{ij}}^n)^{(u)} (\Delta \nu_{ij}^n) \cdot \frac{1}{2} \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} \right], \\ T_{n,1,2,3}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \sum_{u=1}^{d_2} \frac{1}{2} (\Delta W_{e_{ij}}^n)^{(u)} (\Delta \nu_{ij}^n) \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*), \\ T_{n,1,2,4}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \cdot \frac{1}{4} (\Delta \nu_{ij}^n)^2. \end{aligned}$$

Next, we decompose the term $T_{n,2}(Y|X)$ as

$$\begin{aligned} T_{n,2}(Y|X) &:= \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp(g(X, W_{e_i}^n))) - \log(1 + \exp(g(X, W_{e_j}^*))) \right] p_{G_n}(Y|X) \\ &\quad + \sum_{j \in [N^*]: |\mathcal{C}_j| > 1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp(g(X, W_{e_i}^n))) - \log(1 + \exp(g(X, W_{e_j}^*))) \right] p_{G_n}(Y|X) \\ &:= T_{n,2,1}(Y|X) + T_{n,2,2}(Y|X). \end{aligned}$$

Note that we can rewrite the term $T_{n,1,2}(Y|X)$ using the first-order Taylor expansion to the function $\exp(\log(1 + \exp(g(W_{e_i}^n))))$ around the point $W_{e_j}^*$ as

$$T_{n,2,1}(Y|X) = \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \sum_{u=1}^{d_2} (\Delta W_{e_{ij}}^n)^{(u)} \cdot \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} H_n(Y|X; W_{e_j}^*) + R_{n,2,1}(Y|X),$$

where we denote $H_n(Y|X; W_e) = \exp(\log(1 + \exp(g(X, W_e)))p_{G_n}(Y|X)$ and $R_{n,2,1}(Y|X)$ is the Taylor remainder such that $R_{n,2,1}(Y|X)/\mathcal{L}_1(G_n, G_*) \rightarrow 0$ as $n \rightarrow \infty$.

On the other hand, by means of the second-order Taylor expansion, we have

$$T_{n,2,2}(Y|X) = \sum_{j \in [N^*]: |\mathcal{C}_j|>1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\sum_{u=1}^{d_2} (\Delta W_{e_{ij}}^n)^{(u)} \cdot \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} + \sum_{u,v=1}^{d_2} \frac{(\Delta W_{e_{ij}}^n)^{(u)} (\Delta W_{e_{ij}}^n)^{(v)}}{1 + 1_{\{u=v\}}} \cdot \frac{\frac{\partial^2 g}{\partial W_e^{(u)} \partial W_e^{(v)}}(X, W_{e_j}^*) + \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} \right] H_n(Y|X; W_{e_j}^*) + R_{n,2,2}(Y|X),$$

where $R_{n,2,1}(Y|X)$ is the Taylor remainder such that $R_{n,2,2}(Y|X)/\mathcal{L}_1(G_n, G_*) \rightarrow 0$ as $n \rightarrow \infty$.

From the above equation, $[T_{n,1,1}(Y|X) - R_{n,1,1}(Y|X)]$, $[T_{n,1,2}(Y|X) - R_{n,1,2}(Y|X)]$, $[T_{n,2,1}(Y|X) - R_{n,2,1}(Y|X)]$, $[T_{n,2,2}(Y|X) - R_{n,2,2}(Y|X)]$ and $[T_{n,3}(Y|X)]$ can be seen as a combination of elements of the set $\mathcal{S} := \bigcup_{j=1}^N \bigcup_{\rho=0}^5 \mathcal{S}_{\rho,j}$, where we define

$$\begin{aligned} \mathcal{S}_{0,j} &:= \left\{ \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} F_0(Y|X; W_{e_j}^*, \nu_j^*), \frac{\frac{\partial^2 g}{\partial W_e^{(u)} \partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} F_0(Y|X; W_{e_j}^*, \nu_j^*), \right. \\ &\quad \left. \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} F_0(Y|X; W_{e_j}^*, \nu_j^*), F_0(Y|X; W_{e_j}^*, \nu_j^*) : 1 \leq u, v \leq d_2 \right\}, \\ \mathcal{S}_{1,j} &:= \left\{ \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} F_1(Y|X; W_{e_j}^*, \nu_j^*), \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} F_1(Y|X; W_{e_j}^*, \nu_j^*) \right. \\ &\quad \left. \frac{\partial^2 g}{\partial W_e^{(u)} \partial W_e^{(v)}}(X, W_{e_j}^*) F_1(Y|X; W_{e_j}^*, \nu_j^*) : 1 \leq u, v \leq d_2 \right\}, \\ \mathcal{S}_{2,j} &:= \left\{ F_2(Y|X; W_{e_j}^*, \nu_j^*), \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} F_2(Y|X; W_{e_j}^*, \nu_j^*), \right. \\ &\quad \left. \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*) F_2(Y|X; W_{e_j}^*, \nu_j^*) : 1 \leq u, v \leq d_2 \right\}, \\ \mathcal{S}_{3,j} &:= \left\{ \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) F_3(Y|X; W_{e_j}^*, \nu_j^*) : 1 \leq u \leq d_2 \right\}, \\ \mathcal{S}_{4,j} &:= \left\{ F_4(Y|X; W_{e_j}^*, \nu_j^*) \right\}, \\ \mathcal{S}_{5,j} &:= \left\{ \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} H_n(Y|X; W_{e_j}^*, \nu_j^*), \frac{\frac{\partial^2 g}{\partial W_e^{(u)} \partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} H_n(Y|X; W_{e_j}^*, \nu_j^*), \right. \\ &\quad \left. \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} H_n(Y|X; W_{e_j}^*, \nu_j^*), H_n(Y|X; W_{e_j}^*, \nu_j^*) : 1 \leq u, v \leq d_2 \right\}. \end{aligned}$$

Step 2: Non-vanishing coefficients. In this step, we will show that at least one among the coefficients in the representations of $[T_{n,1,1}(Y|X) - R_{n,1,1}(Y|X)]/\mathcal{L}_1(G_n, G_*)$, $[T_{n,1,2}(Y|X) - R_{n,1,2}(Y|X)]/\mathcal{L}_1(G_n, G_*)$, $[T_{n,2,1}(Y|X) - R_{n,2,1}(Y|X)]/\mathcal{L}_1(G_n, G_*)$, $[T_{n,2,2}(Y|X) - R_{n,2,2}(Y|X)]/\mathcal{L}_1(G_n, G_*)$ and $[T_{n,3}(Y|X)]/\mathcal{L}_1(G_n, G_*)$ does not approach zero when n goes to infinity. Assume by contrary that all of them vanish as $n \rightarrow \infty$. Then, by considering the coefficients of the term

- $F_0(Y|X; W_{e_j}^*, \nu_j^*)$ for $j \in [N^*]$, we have

$$\frac{1}{\mathcal{L}_1(G_n, G_*)} \cdot \sum_{j=1}^{N^*} \left| \sum_{i \in \mathcal{C}_j} \exp(c_i^n) - \exp(c_j^*) \right| \rightarrow 0.$$

- $\frac{\frac{\partial g}{\partial W_{e_j}^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} F_0(Y|X; W_{e_j}^*, \nu_j^*)$ for $j \in [N^*] : |\mathcal{C}_j| = 1$, we have

$$\frac{1}{\mathcal{L}_1(G_n, G_*)} \cdot \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \|\Delta W_{e_{ij}}^n\|_1 \rightarrow 0.$$

Due to the equivalence between the ℓ_1 -norm and the ℓ_2 -norm, we obtain

$$\frac{1}{\mathcal{L}_1(G_n, G_*)} \cdot \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \|\Delta W_{e_{ij}}^n\| \rightarrow 0.$$

- $F_2(Y|X; W_{e_j}^*, \nu_j^*)$ for $j \in [N^*] : |\mathcal{C}_j| = 1$, we have

$$\frac{1}{\mathcal{L}_1(G_n, G_*)} \cdot \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) |\Delta \nu_{ij}^n| \rightarrow 0.$$

- $\frac{\frac{\partial g}{\partial W_{e_j}^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_{e_j}^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} F_0(Y|X; W_{e_j}^*, \nu_j^*)$ for $j \in [N^*] : |\mathcal{C}_j| > 1$, we have

$$\frac{1}{\mathcal{L}_1(G_n, G_*)} \cdot \sum_{j \in [N^*]: |\mathcal{C}_j|>1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \|\Delta W_{e_{ij}}^n\|^2 \rightarrow 0.$$

- $F_4(Y|X; W_{e_j}^*, \nu_j^*)$ for $j \in [N^*] : |\mathcal{C}_j| > 1$, we have

$$\frac{1}{\mathcal{L}_1(G_n, G_*)} \cdot \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) |\Delta \nu_{ij}^n|^2 \rightarrow 0.$$

By taking the sum of the above limits, we obtain $1 = \frac{\mathcal{L}_1(G_n, G_*)}{\mathcal{L}_1(G_n, G_*)} \rightarrow 0$ as $n \rightarrow \infty$, which is a contradiction. Thus, not all the coefficients in the representations of $[T_{n,1,1}(Y|X) - R_{n,1,1}(Y|X)]/\mathcal{L}_1(G_n, G_*)$, $[T_{n,1,2}(Y|X) - R_{n,1,2}(Y|X)]/\mathcal{L}_1(G_n, G_*)$, $[T_{n,2,1}(Y|X) - R_{n,2,1}(Y|X)]/\mathcal{L}_1(G_n, G_*)$, $[T_{n,2,2}(Y|X) - R_{n,2,2}(Y|X)]/\mathcal{L}_1(G_n, G_*)$ and $[T_{n,3}(Y|X)]/\mathcal{L}_1(G_n, G_*)$ converge to zero as $n \rightarrow \infty$.

Stage 3 - Fatou's argument: In this stage, we use the Fatou's lemma to show a contradiction to the result of Step 2. For that purpose, let us denote m_n as the maximum of the absolute values of the coefficients in the representations of $[T_{n,1,1}(Y|X) - R_{n,1,1}(Y|X)]/\mathcal{L}_1(G_n, G_*)$, $[T_{n,1,2}(Y|X) - R_{n,1,2}(Y|X)]/\mathcal{L}_1(G_n, G_*)$, $[T_{n,2,1}(Y|X) - R_{n,2,1}(Y|X)]/\mathcal{L}_1(G_n, G_*)$, $[T_{n,2,2}(Y|X) - R_{n,2,2}(Y|X)]/\mathcal{L}_1(G_n, G_*)$ and $[T_{n,3}(Y|X)]/\mathcal{L}_1(G_n, G_*)$. It follows from the result of Step 2 that $1/m_n \not\rightarrow \infty$ as $n \rightarrow \infty$. In addition, we also denote

$$\begin{aligned} \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta W_{e_{ij}}^n)^{(u)}}{m_n \mathcal{L}_1(G_n, G_*)} &\rightarrow \alpha_{1,j}^{(u)}, & \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta \nu_{ij}^n)}{m_n \mathcal{L}_1(G_n, G_*)} &\rightarrow \beta_{1,j}, \\ \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta W_{e_{ij}}^n)^{(u)} (\Delta W_{e_{ij}}^n)^{(v)}}{m_n \mathcal{L}_1(G_n, G_*)} &\rightarrow \alpha_{2,j}^{(uv)}, & \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta \nu_{ij}^n)^2}{m_n \mathcal{L}_1(G_n, G_*)} &\rightarrow \beta_{2,j}, \\ \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta W_{e_{ij}}^n)^{(u)} (\Delta \nu_{ij}^n)}{m_n \mathcal{L}_1(G_n, G_*)} &\rightarrow \gamma_j^{(u)}, & \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) - \exp(c_j^*)}{m_n \mathcal{L}_1(G_n, G_*)} &\rightarrow \xi_j, \end{aligned}$$

as $n \rightarrow \infty$ for any $j \in [N^*]$ and $u, v \in [d_2]$ with a note that at least one among $\alpha_{1,j}^{(u)}, \beta_{1,j}, \alpha_{2,j}^{(uv)}, \beta_{2,j}, \gamma_j^{(u)}$ and ξ_j is non-zero.

By applying the Fatou's lemma, we have

$$0 = \lim_{n \rightarrow \infty} \frac{\mathbb{E}_X[V(p_{G_n}(\cdot|X), p_{G_*}(\cdot|X))]}{m_n \mathcal{L}_1(G_n, G_*)} = \frac{1}{2} \int \liminf_{n \rightarrow \infty} \frac{|p_{G_n}(Y|X) - p_{G_*}(Y|X)|}{m_n \mathcal{L}_1(G_n, G_*)} d(X, Y),$$

which implies that $[p_{G_n}(Y|X) - p_{G_*}(Y|X)]/[m_n \mathcal{L}_1(G_n, G_*)] \rightarrow 0$ as $n \rightarrow \infty$ for almost surely (X, Y) . Since the term $\sum_{j=1}^{N^*} \exp(\log(1 + \exp(g(x, W_{e_j}^*)))$ is bounded, we also have $T_n(Y|X)/[m_n \mathcal{L}_1(G_n, G_*)] \rightarrow 0$ as $n \rightarrow \infty$. Then, it follows that

$$0 = \lim_{n \rightarrow \infty} \frac{T_{n,1,1}(Y|X) + T_{n,1,2}(Y|X)}{m_n \mathcal{L}_1(G_n, G_*)} - \lim_{n \rightarrow \infty} \frac{T_{n,2,1}(Y|X) + T_{n,2,2}(Y|X)}{m_n \mathcal{L}_1(G_n, G_*)} + \lim_{n \rightarrow \infty} \frac{T_{n,3}(Y|X)}{m_n \mathcal{L}_1(G_n, G_*)}, \quad (20)$$

for almost surely $(X, Y) \in \mathcal{X} \times \mathcal{Y}$, where we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{T_{n,1,1}(Y|X)}{m_n \mathcal{L}_1(G_n, G_*)} &:= \sum_{j \in [N^*]; |\mathcal{C}_j|=1} \left[\sum_{u=1}^{d_2} \alpha_{1,j}^{(u)} \frac{\frac{\partial g}{\partial W_{e_j}^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} F_{0,j}(Y|X) \right. \\ &\quad \left. + \sum_{u=1}^{d_2} \alpha_{1,j}^{(u)} \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) F_{1,j}(Y|X) + \frac{1}{2} \beta_{1,j} F_{2,j}(Y|X) \right], \\ \lim_{n \rightarrow \infty} \frac{T_{n,1,2}(Y|X)}{m_n \mathcal{L}_1(G_n, G_*)} &:= \sum_{j \in [N^*]; |\mathcal{C}_j|>1} \left[\left(\sum_{u=1}^{d_2} \alpha_{1,j}^{(u)} \frac{\frac{\partial g}{\partial W_{e_j}^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} \right. \right. \\ &\quad \left. \left. + \sum_{u,v=1}^{d_2} \frac{\alpha_{2,j}^{(uv)}}{1 + 1_{\{u=v\}}} \cdot \frac{\frac{\partial^2 g}{\partial W_e^{(u)} \partial W_e^{(v)}}(X, W_{e_j}^*) + \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} \right) F_{0,j}(Y|X) \right. \\ &\quad \left. + \left(\sum_{u=1}^{d_2} \alpha_{1,j}^{(u)} \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) + \sum_{u,v=1}^{d_2} \frac{\alpha_{2,j}^{(uv)}}{1 + 1_{\{u=v\}}} \left(\frac{2 \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} \right. \right. \right. \\ &\quad \left. \left. + \frac{\partial^2 g}{\partial W_e^{(u)} \partial W_e^{(v)}}(X, W_{e_j}^*) \right) \right) F_{1,j}(Y|X) + \left(\frac{1}{2} \beta_{1,j} + \sum_{u=1}^{d_2} \gamma_j^{(u)} \cdot \frac{1}{2} \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} \right. \\ &\quad \left. + \sum_{u,v=1}^{d_2} \frac{\alpha_{2,j}^{(uv)}}{1 + 1_{\{u=v\}}} \cdot \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*) \right) F_{2,j}(Y|X) \\ &\quad \left. + \sum_{u=1}^{d_2} \frac{1}{2} \gamma_j^{(u)} \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) F_{3,j}(Y|X) + \frac{1}{4} \beta_{2,j} F_{4,j}(Y|X) \right], \end{aligned}$$

and

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{T_{n,2,1}(Y|X)}{m_n \mathcal{L}_1(G_n, G_*)} &:= \sum_{j \in [N^*]; |\mathcal{C}_j|=1} \sum_{u=1}^{d_2} \alpha_{1,j}^{(u)} \cdot \frac{\frac{\partial g}{\partial W_{e_j}^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} H_j(Y|X), \\ \lim_{n \rightarrow \infty} \frac{T_{n,2,2}(Y|X)}{m_n \mathcal{L}_1(G_n, G_*)} &:= \sum_{j \in [N^*]; |\mathcal{C}_j|>1} \left[\sum_{u=1}^{d_2} \alpha_{1,j}^{(u)} \cdot \frac{\frac{\partial g}{\partial W_{e_j}^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} \right. \\ &\quad \left. + \sum_{u,v=1}^{d_2} \frac{\alpha_{2,j}^{(uv)}}{1 + 1_{\{u=v\}}} \cdot \frac{\frac{\partial^2 g}{\partial W_e^{(u)} \partial W_e^{(v)}}(X, W_{e_j}^*) + \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} \right] H_j(Y|X), \end{aligned}$$

and

$$\lim_{n \rightarrow \infty} \frac{T_{n,3}(Y|X)}{m_n \mathcal{L}_1(G_n, G_*)} := \sum_{j=1}^{N^*} \xi_j [F_{0,j}(Y|X) - H_j(Y|X)].$$

It is worth noting that for almost every X , the set

$$\left\{ F_{\rho,j}(Y|X), H_j(Y|X) : 0 \leq \rho \leq 4, j \in [N^*] \right\}$$

is linearly independent w.r.t Y . Therefore, it follows that the coefficients of those terms in the limit in equation (20) become zero.

For $j \in [N^*]$ such that $|\mathcal{C}_j| = 1$, by considering the coefficients of

- $F_{0,j}(Y|X)$, we have $\xi_j + \sum_{u=1}^{d_2} \alpha_{1,j}^{(u)} \cdot \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} = 0$ for almost surely X . Since the expert function g is strongly identifiable, we deduce $\xi_j = \alpha_{1,j}^{(u)} = 0$ for all $u \in [d_2]$;
- $F_{2,j}(Y|X)$, we have $\beta_{1,j} = 0$.

For $j \in [N^*]$ such that $|\mathcal{C}_j| > 1$, by considering the coefficients of

- $F_{0,j}(Y|X)$, we have

$$\xi_j + \sum_{u=1}^{d_2} \alpha_{1,j}^{(u)} \frac{\frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} + \sum_{u,v=1}^{d_2} \frac{\alpha_{2,j}^{(uv)}}{1 + 1_{\{u=v\}}} \cdot \frac{\frac{\partial^2 g}{\partial W_e^{(u)} \partial W_e^{(v)}}(X, W_{e_j}^*) + \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) \frac{\partial g}{\partial W_e^{(v)}}(X, W_{e_j}^*)}{1 + \exp(-g(X, W_{e_j}^*))} = 0$$

for almost surely X . Since the expert function g is strongly identifiable, we deduce $\xi_j = \alpha_{1,j}^{(u)} = \alpha_{2,j}^{(uv)} = 0$ for all $u, v \in [d_2]$;

- $F_{3,j}(Y|X)$, we have $\sum_{u=1}^{d_2} \frac{1}{2} \gamma_j^{(u)} \frac{\partial g}{\partial W_e^{(u)}}(X, W_{e_j}^*) = 0$ for almost surely X . Since the expert function g is strongly identifiable, we deduce $\gamma_j^{(u)} = 0$ for all $u \in [d_2]$;
- $F_{4,j}(Y|X)$, we have $\beta_{2,j} = 0$.

Putting the above results together, we have (ii) $\xi_j = \alpha_{1,j}^{(u)} = \beta_{1,j} = \alpha_{2,j}^{(uv)} = \beta_{2,j} = \gamma_j^{(u)} = 0$ for all $j \in [N^*]$ and $u, v \in [d_2]$. This contradicts to the fact that at least one among them is non-zero. Consequently, we achieve the local part in equation (17).

Global part: Now, it suffices to demonstrate that

$$\inf_{G \in \mathcal{G}_N(\Theta) : \mathcal{L}_1(G, G_*) > \varepsilon'} \frac{\mathbb{E}_X[V(p_G(\cdot|X), p_{G_*}(\cdot|X))]}{\mathcal{L}_1(G, G_*)} > 0,$$

for some positive constant ε' . Given the above result, it is sufficient to derive the global part in equation (18), that is,

$$\inf_{G \in \mathcal{G}_N(\Theta) : \mathcal{L}_1(G, G_*) > \varepsilon'} \mathbb{E}_X[V(p_G(\cdot|X), p_{G_*}(\cdot|X))] / \mathcal{L}_1(G, G_*) > 0.$$

Assume by contrary that the global part does not hold true, then we can find a sequence $\tilde{G}_n \in \mathcal{G}_N(\Theta)$ such that $\mathcal{L}_1(\tilde{G}_n, G_*) > \varepsilon'$ and $\mathbb{E}_X[V(p_{\tilde{G}_n}(\cdot|X), p_{G_*}(\cdot|X))] \rightarrow 0$ as $n \rightarrow \infty$. Since Θ is a compact set, we are able to replace \tilde{G}_n with its subsequence which converges to some mixing measure $\tilde{G} \in \mathcal{G}_N(\Theta)$. Recall that $\mathcal{L}_1(\tilde{G}_n, G_*) > \varepsilon'$, then we also get that $\mathcal{L}_1(\tilde{G}, G_*) > \varepsilon'$.

On the other hand, by means of the Fatou's lemma, we have

$$0 = \lim_{n \rightarrow \infty} \mathbb{E}_X[2V(p_{\tilde{G}_n}(\cdot|X), p_{G_*}(\cdot|X))] \geq \int \liminf_{n \rightarrow \infty} |p_{\tilde{G}_n}(Y|X) - p_{G_*}(Y|X)| d(X, Y),$$

which follows that $p_{\tilde{G}}(Y|X) - p_{G_*}(Y|X) = 0$ for almost surely (X, Y) . Thus, we achieve that $\tilde{G} \equiv G_*$, or equivalently $\mathcal{L}_1(\tilde{G}, G_*) = 0$. This contradicts to the fact that $\mathcal{L}_1(\tilde{G}, G_*) > \varepsilon' > 0$.

Hence, we reach the conclusion in equation (18), and the proof is completed.

2052 **K.2 PROOF OF THEOREM J.1**

2053 As in Appendix K.1, we also start with establishing the local part

2054
$$\lim_{\varepsilon \rightarrow 0} \inf_{G \in \mathcal{G}_N(\Theta): \mathcal{L}_2(G, G_*) \leq \varepsilon} \frac{\mathbb{E}_X[V(p_G(\cdot|X), p_{G_*}(\cdot|X))]}{\mathcal{L}_2(G, G_*)} > 0. \quad (21)$$

2055 Assume by contrary that the local part is not true, then we can find a sequence of mixing measures
 2056 (G_n) given by $G_n := \sum_{i=1}^N \exp(c_i^n) \delta_{(a_i^n, b_i^n, \nu_i^n)} \in \mathcal{G}_N(\Theta)$ such that $\mathcal{L}_2(G_n, G_*) \rightarrow 0$ and

2057
$$\mathbb{E}_X[V(p_{G_n}(\cdot|X), p_{G_*}(\cdot|X))]/\mathcal{L}_2(G_n, G_*) \rightarrow 0,$$

2058 as $n \rightarrow \infty$. Recall that the Voronoi loss $\mathcal{L}_2(G_n, G_*)$ is given by

2059
$$\begin{aligned} \mathcal{L}_2(G_n, G_*) := & \sum_{j=1}^{N^*} \left| \sum_{i \in \mathcal{C}_j} \exp(c_i^n) - \exp(c_j^*) \right| + \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\|W_{e_i}^n - W_{e_j}^*\| + |\nu_i^n - \nu_j^*| \right] \\ & + \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\|W_{e_i}^n - W_{e_j}^*\|^2 + |\nu_i^n - \nu_j^*|^2 \right]. \end{aligned} \quad (22)$$

2060 Since $\mathcal{L}_2(G_n, G_*) \rightarrow 0$ as $n \rightarrow \infty$, we obtain $(a_i^n, b_i^n, \nu_i^n) \rightarrow (a_j^*, b_j^*, \nu_j^*)$ for all $j \in [N^*]$ and
 2061 $i \in \mathcal{C}_j$.

2062 Next, we divide the rest of this proof into three main steps.

2063 **Step 1: Taylor expansion.** In this step, we aim to decompose the term $T_n(Y|X) :=$
 2064 $\left[\sum_{j=1}^{N^*} \exp(\log(1 + \exp((a_j^*)^\top X + b_j^*))) \right] \cdot [p_{G_n}(Y|X) - p_{G_*}(Y|X)]$ can be decomposed as

2065
$$\begin{aligned} T_n(Y|X) = & \sum_{j=1}^{N^*} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp((a_i^n)^\top X + b_i^n))) f(Y|(a_i^n)^\top X + b_i^n, \nu_i^n) \right. \\ & \left. - \exp(\log(1 + \exp((a_j^*)^\top X + b_j^*))) f(Y|(a_j^*)^\top X + b_j^*, \nu_j^*) \right] \\ & - \sum_{j=1}^{N^*} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp((a_i^n)^\top X + b_i^n))) - \exp(\log(1 + \exp((a_j^*)^\top X + b_j^*))) \right] p_{G_n}(Y|X) \\ & + \sum_{j=1}^{N^*} \left[\sum_{i \in \mathcal{C}_j} \exp(c_i^n) - \exp(c_j^*) \right] \cdot \exp(\log(1 + \exp((a_j^*)^\top X + b_j^*))) [f(Y|(a_j^*)^\top X + b_j^*, \nu_j^*) - p_{G_n}(Y|X)] \\ := & T_{n,1}(Y|X) - T_{n,2}(Y|X) + T_{n,3}(Y|X). \end{aligned}$$

2066 Next, we continue to decompose the term $T_{n,1}(Y|X)$ as

2067
$$\begin{aligned} T_{n,1}(Y|X) = & \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp((a_i^n)^\top X + b_i^n))) f(Y|(a_i^n)^\top X + b_i^n, \nu_i^n) \right. \\ & \left. - \exp(\log(1 + \exp((a_j^*)^\top X + b_j^*))) f(Y|(a_j^*)^\top X + b_j^*, \nu_j^*) \right] \\ & + \sum_{j \in [N^*]: |\mathcal{C}_j|>1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp((a_i^n)^\top X + b_i^n))) f(Y|(a_i^n)^\top X + b_i^n, \nu_i^n) \right. \\ & \left. - \exp(\log(1 + \exp((a_j^*)^\top X + b_j^*))) f(Y|(a_j^*)^\top X + b_j^*, \nu_j^*) \right] \\ := & T_{n,1,1}(Y|X) + T_{n,1,2}(Y|X). \end{aligned}$$

2068 Let us denote $F_\rho(Y|X; a, b, \nu) := \exp(\log(1 + \exp(a^\top X + b))) \frac{\partial^\rho f}{\partial g^\rho}(Y|a^\top X + b, \nu)$. By applying
 2069 the first-order Taylor expansion to the function $F_0(Y|X; a, b, \nu)$ around the point (a_j^*, b_j^*, ν_j^*) , we
 2070 rewrite the term $T_{n,1,1}(Y|X)$ as

2071
$$T_{n,1,1}(Y|X) = \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{\rho=0}^2 T_{n,1,1,\rho}^{(j)}(X) F_\rho(Y; X, a_j^*, b_j^*, \nu_j^*) + R_{n,1,1}(Y|X),$$

where $R_{n,1,1}(Y|X)$ is the Taylor remainder such that $R_{n,1,1}(Y|X)/\mathcal{L}_2(G_n, G_*) \rightarrow 0$ as $n \rightarrow \infty$, and

$$\begin{aligned} T_{n,1,1,0}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \cdot \frac{\sum_{u=1}^d (\Delta a_{ij}^n)^{(u)} X^{(u)} + (\Delta b_{ij}^n)}{1 + \exp(-(a_j^*)^\top X - b_j^*)}, \\ T_{n,1,1,1}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\sum_{u=1}^d (\Delta a_{ij}^n)^{(u)} X^{(u)} + (\Delta b_{ij}^n) \right], \\ T_{n,1,1,2}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \frac{1}{2} \exp(c_i^n) (\Delta \nu_{ij}^n), \end{aligned}$$

in which $\Delta a_{ij}^n := a_i^n - a_j^*$, $\Delta b_{ij}^n := b_i^n - b_j^*$ and $\Delta \nu_{ij}^n := \nu_i^n - \nu_j^*$.

Meanwhile, by means of the second-order Taylor expansion, the term $T_{n,1,2}(Y|X)$ can be represented as

$$T_{n,1,2}(Y|X) = \sum_{j \in [N^*]: |\mathcal{C}_j| > 1} \sum_{\rho=0}^4 T_{n,1,2,\rho}^{(j)}(X) F_\rho(Y; X, a_j^*, b_j^*, \nu_j^*) + R_{n,1,2}(Y|X),$$

where $R_{n,1,2}(Y|X)$ is the Taylor remainder such that $R_{n,1,2}(Y|X)/\mathcal{L}_2(G_n, G_*) \rightarrow 0$ as $n \rightarrow \infty$, and

$$\begin{aligned} T_{n,1,2,0}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\frac{\sum_{u=1}^d (\Delta a_{ij}^n)^{(u)} X^{(u)} + (\Delta b_{ij}^n)}{1 + \exp(-(a_j^*)^\top X - b_j^*)} + \frac{\sum_{u,v=1}^d \frac{(\Delta a_{ij}^n)^{(u)} (\Delta a_{ij}^n)^{(v)} X^{(u)} X^{(v)}}{1 + \mathbb{1}_{\{u=v\}}} X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right. \\ &\quad \left. + \frac{\sum_{u=1}^d (\Delta a_{ij}^n)^{(u)} (\Delta b_{ij}^n) X^{(u)} + \frac{1}{2} (\Delta b_{ij}^n)^2}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right], \\ T_{n,1,2,1}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\sum_{u=1}^d (\Delta a_{ij}^n)^{(u)} X^{(u)} + (\Delta b_{ij}^n) + \frac{2 \sum_{u,v=1}^d \frac{(\Delta a_{ij}^n)^{(u)} (\Delta a_{ij}^n)^{(v)} X^{(u)} X^{(v)}}{1 + \mathbb{1}_{\{u=v\}}}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right. \\ &\quad \left. + \frac{(\Delta b_{ij}^n)^2 + 2 \sum_{u=1}^d (\Delta a_{ij}^n)^{(u)} (\Delta b_{ij}^n) X^{(u)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right], \\ T_{n,1,2,2}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\frac{1}{2} (\Delta \nu_{ij}^n) + \sum_{u,v=1}^d \frac{(\Delta a_{ij}^n)^{(u)} (\Delta a_{ij}^n)^{(v)} X^{(u)} X^{(v)}}{1 + \mathbb{1}_{\{u=v\}}} + \frac{1}{2} (\Delta b_{ij}^n)^2 \right. \\ &\quad \left. + \sum_{u=1}^d (\Delta a_{ij}^n)^{(u)} (\Delta b_{ij}^n) X^{(u)} + \frac{1}{2} \cdot \frac{\sum_{u=1}^d (\Delta a_{ij}^n)^{(u)} (\Delta \nu_{ij}^n) X^{(u)} + (\Delta b_{ij}^n) (\Delta \nu_{ij}^n)}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right], \\ T_{n,1,2,3}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\sum_{u=1}^d \frac{1}{2} (\Delta a_{ij}^n)^{(u)} (\Delta \nu_{ij}^n) X^{(u)} + \frac{1}{2} (\Delta b_{ij}^n) (\Delta \nu_{ij}^n) \right], \\ T_{n,1,2,4}^{(j)}(X) &:= \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \cdot \frac{1}{4} (\Delta \nu_{ij}^n)^2. \end{aligned}$$

Next, we decompose the term $T_{n,2}(Y|X)$ as

$$\begin{aligned} T_{n,2}(Y|X) &:= \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp((a_i^n)^\top X + b_i^n))) - \exp(\log(1 + \exp((a_j^*)^\top X + b_j^*))) \right] p_{G_n}(Y|X) \\ &\quad + \sum_{j \in [N^*]: |\mathcal{C}_j| > 1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\exp(\log(1 + \exp((a_i^n)^\top X + b_i^n))) - \exp(\log(1 + \exp((a_j^*)^\top X + b_j^*))) \right] p_{G_n}(Y|X) \\ &:= T_{n,2,1}(Y|X) + T_{n,2,2}(Y|X). \end{aligned}$$

Note that we can rewrite the term $T_{n,1,2}(Y|X)$ using the first-order Taylor expansion to the function $\exp(\log(1 + \exp((a_i^n)^\top X + b_i^n)))$ around the point (a_j^*, b_j^*) as

$$T_{n,2,1}(Y|X) = \sum_{j \in [N^*]: |\mathcal{C}_j| = 1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \cdot \frac{\sum_{u=1}^d (\Delta a_{ij}^n)^{(u)} X^{(u)} + (\Delta b_{ij}^n)}{1 + \exp(-(a_j^*)^\top X - b_j^*)} H_n(Y|X; a_j^*, b_j^*) + R_{n,2,1}(Y|X),$$

where we denote $H_n(Y|X; a, b) = \exp(\log(1 + \exp(a^\top X + b))) p_{G_n}(Y|X)$ and $R_{n,2,1}(Y|X)$ is the Taylor remainder such that $R_{n,2,1}(Y|X)/\mathcal{L}_2(G_n, G_*) \rightarrow 0$ as $n \rightarrow \infty$.

On the other hand, by means of the second-order Taylor expansion, we have

$$T_{n,2,2}(Y|X) = \sum_{j \in [N^*]: |\mathcal{C}_j| > 1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \left[\frac{\sum_{u=1}^d (\Delta a_{ij}^n)^{(u)} X^{(u)} + (\Delta b_{ij}^n)}{1 + \exp(-(a_j^*)^\top X - b_j^*)} + \frac{\sum_{u,v=1}^d \frac{(\Delta a_{ij}^n)^{(u)} (\Delta a_{ij}^n)^{(v)}}{1 + 1_{\{u=v\}}} X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} + \frac{\sum_{u=1}^d (\Delta a_{ij}^n)^{(u)} (\Delta b_{ij}^n) X^{(u)} + \frac{1}{2} (\Delta b_{ij}^n)^2}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right] H_n(Y|X; W_{e_j^*}) + R_{n,2,2}(Y|X),$$

where $R_{n,2,1}(Y|X)$ is the Taylor remainder such that $R_{n,2,2}(Y|X)/\mathcal{L}_2(G_n, G_*) \rightarrow 0$ as $n \rightarrow \infty$.

From the above equation, $[T_{n,1,1}(Y|X) - R_{n,1,1}(Y|X)]$, $[T_{n,1,2}(Y|X) - R_{n,1,2}(Y|X)]$, $[T_{n,2,1}(Y|X) - R_{n,2,1}(Y|X)]$, $[T_{n,2,2}(Y|X) - R_{n,2,2}(Y|X)]$ and $[T_{n,3}(Y|X)]$ can be seen as a combination of elements of the set $\mathcal{S} := \bigcup_{j=1}^N \bigcup_{\rho=0}^5 \mathcal{S}_{\rho,j}$, where we define

$$\begin{aligned} \mathcal{S}_{0,j} &:= \left\{ \frac{X^{(u)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{0,j}(Y|X), \frac{X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{0,j}(Y|X), \right. \\ &\quad \left. \frac{1}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{0,j}(Y|X), F_{0,j}(Y|X) : 1 \leq u, v \leq d \right\}, \\ \mathcal{S}_{1,j} &:= \left\{ F_{1,j}(Y|X), X^{(u)} F_{1,j}(Y|X), \frac{X^{(u)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{1,j}(Y|X), \right. \\ &\quad \left. \frac{X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{1,j}(Y|X), \frac{1}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{1,j}(Y|X) : 1 \leq u, v \leq d \right\}, \\ \mathcal{S}_{2,j} &:= \left\{ F_{2,j}(Y|X), X^{(u)} F_{2,j}(Y|X), X^{(u)} X^{(v)} F_{2,j}(Y|X), \right. \\ &\quad \left. \frac{X^{(u)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{2,j}(Y|X), \frac{1}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{2,j}(Y|X) : 1 \leq u, v \leq d_2 \right\}, \\ \mathcal{S}_{3,j} &:= \left\{ F_{3,j}(Y|X), X^{(u)} F_{3,j}(Y|X) : 1 \leq u \leq d \right\}, \\ \mathcal{S}_{4,j} &:= \left\{ F_{4,j}(Y|X) \right\}, \\ \mathcal{S}_{5,j} &:= \left\{ \frac{X^{(u)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} H_{n,j}(Y|X), \frac{X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} H_{n,j}(Y|X), \right. \\ &\quad \left. \frac{1}{1 + \exp(-(a_j^*)^\top X - b_j^*)} H_{n,j}(Y|X), H_{n,j}(Y|X) : 1 \leq u, v \leq d \right\}. \end{aligned}$$

Step 2: Non-vanishing coefficients. In this step, we will show that at least one among the coefficients in the representations of $[T_{n,1,1}(Y|X) - R_{n,1,1}(Y|X)]/\mathcal{L}_2(G_n, G_*)$,

2214 $[T_{n,1,2}(Y|X) - R_{n,1,2}(Y|X)]/\mathcal{L}_2(G_n, G_*)$, $[T_{n,2,1}(Y|X) - R_{n,2,1}(Y|X)]/\mathcal{L}_2(G_n, G_*)$,
 2215 $[T_{n,2,2}(Y|X) - R_{n,2,2}(Y|X)]/\mathcal{L}_2(G_n, G_*)$ and $[T_{n,3}(Y|X)]/\mathcal{L}_2(G_n, G_*)$ does not approach zero
 2216 when n goes to infinity. Assume by contrary that all of them vanish as $n \rightarrow \infty$. Then, by considering
 2217 the coefficients of the term

- $F_{0,j}(Y|X)$ for $j \in [N^*]$, we have

$$\frac{1}{\mathcal{L}_2(G_n, G_*)} \cdot \sum_{j=1}^{N^*} \left| \sum_{i \in \mathcal{C}_j} \exp(c_i^n) - \exp(c_j^*) \right| \rightarrow 0.$$

- $\frac{X^{(u)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{0,j}(Y|X)$ for $j \in [N^*] : |\mathcal{C}_j| = 1$, we have

$$\frac{1}{\mathcal{L}_2(G_n, G_*)} \cdot \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \|\Delta a_{ij}^n\| \rightarrow 0.$$

- $\frac{1}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{0,j}(Y|X)$ for $j \in [N^*] : |\mathcal{C}_j| = 1$, we have

$$\frac{1}{\mathcal{L}_2(G_n, G_*)} \cdot \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \|\Delta b_{ij}^n\| \rightarrow 0.$$

- $F_{2,j}(Y|X)$ for $j \in [N^*] : |\mathcal{C}_j| = 1$, we have

$$\frac{1}{\mathcal{L}_2(G_n, G_*)} \cdot \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) |\Delta \nu_{ij}^n| \rightarrow 0.$$

- $\frac{X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{0,j}(Y|X)$ for $j \in [N^*] : |\mathcal{C}_j| > 1$, we have

$$\frac{1}{\mathcal{L}_2(G_n, G_*)} \cdot \sum_{j \in [N^*]: |\mathcal{C}_j|>1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) \|\Delta a_j^n\|^2 \rightarrow 0.$$

- $\frac{1}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{1,j}(Y|X)$ for $j \in [N^*] : |\mathcal{C}_j| > 1$, we have

$$\frac{1}{\mathcal{L}_2(G_n, G_*)} \cdot \sum_{j \in [N^*]: |\mathcal{C}_j|>1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) |\Delta b_j^n|^2 \rightarrow 0.$$

- $F_{4,j}(Y|X)$ for $j \in [N^*] : |\mathcal{C}_j| > 1$, we have

$$\frac{1}{\mathcal{L}_2(G_n, G_*)} \cdot \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \sum_{i \in \mathcal{C}_j} \exp(c_i^n) |\Delta \nu_{ij}^n|^2 \rightarrow 0.$$

2258 By taking the sum of the above limits, we obtain $1 = \frac{\mathcal{L}_2(G_n, G_*)}{\mathcal{L}_2(G_n, G_*)} \rightarrow 0$ as
 2259 $n \rightarrow \infty$, which is a contradiction. Thus, not all the coefficients in the representa-
 2260 tions of $[T_{n,1,1}(Y|X) - R_{n,1,1}(Y|X)]/\mathcal{L}_2(G_n, G_*)$, $[T_{n,1,2}(Y|X) - R_{n,1,2}(Y|X)]/\mathcal{L}_2(G_n, G_*)$,
 2261 $[T_{n,2,1}(Y|X) - R_{n,2,1}(Y|X)]/\mathcal{L}_2(G_n, G_*)$, $[T_{n,2,2}(Y|X) - R_{n,2,2}(Y|X)]/\mathcal{L}_2(G_n, G_*)$ and
 2262 $[T_{n,3}(Y|X)]/\mathcal{L}_2(G_n, G_*)$ converge to zero as $n \rightarrow \infty$.

2263 **Stage 3 - Fatou's argument:** In this stage, we use the Fatou's lemma to show
 2264 a contradiction to the result of Step 2. For that purpose, let us denote m_n
 2265 as the maximum of the absolute values of the coefficients in the representations
 2266 of $[T_{n,1,1}(Y|X) - R_{n,1,1}(Y|X)]/\mathcal{L}_2(G_n, G_*)$, $[T_{n,1,2}(Y|X) - R_{n,1,2}(Y|X)]/\mathcal{L}_2(G_n, G_*)$,
 2267 $[T_{n,2,1}(Y|X) - R_{n,2,1}(Y|X)]/\mathcal{L}_2(G_n, G_*)$, $[T_{n,2,2}(Y|X) - R_{n,2,2}(Y|X)]/\mathcal{L}_2(G_n, G_*)$ and
 $[T_{n,3}(Y|X)]/\mathcal{L}_2(G_n, G_*)$. It follows from the result of Step 2 that $1/m_n \not\rightarrow \infty$ as $n \rightarrow \infty$.

In addition, we also denote

$$\begin{aligned}
& \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta a_{ij}^n)^{(u)}}{m_n \mathcal{L}_2(G_n, G_*)} \rightarrow \alpha_{1,j}^{(u)}, & \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta \nu_{ij}^n)}{m_n \mathcal{L}_2(G_n, G_*)} \rightarrow \beta_{1,j}, \\
& \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta a_{ij}^n)^{(u)} (\Delta a_{ij}^n)^{(v)}}{m_n \mathcal{L}_2(G_n, G_*)} \rightarrow \alpha_{2,j}^{(uv)}, & \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta \nu_{ij}^n)^2}{m_n \mathcal{L}_2(G_n, G_*)} \rightarrow \beta_{2,j}, \\
& \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta b_{ij}^n)}{m_n \mathcal{L}_2(G_n, G_*)} \rightarrow \phi_{1,j}^{(u)}, & \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta b_{ij}^n)^2}{m_n \mathcal{L}_2(G_n, G_*)} \rightarrow \phi_{2,j}, \\
& \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta a_{ij}^n)^{(u)} (\Delta \nu_{ij}^n)}{m_n \mathcal{L}_2(G_n, G_*)} \rightarrow \gamma_{1,j}^{(u)}, & \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta a_{ij}^n)^{(u)} (\Delta b_{ij}^n)}{m_n \mathcal{L}_2(G_n, G_*)} \rightarrow \gamma_{2,j}^{(u)}, \\
& \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) (\Delta b_{ij}^n) (\Delta \nu_{ij}^n)}{m_n \mathcal{L}_2(G_n, G_*)} \rightarrow \gamma_{3,j}, & \frac{\sum_{i \in \mathcal{C}_j} \exp(c_i^n) - \exp(c_j^*)}{m_n \mathcal{L}_2(G_n, G_*)} \rightarrow \xi_j,
\end{aligned}$$

as $n \rightarrow \infty$ for any $j \in [N^*]$ and $u, v \in [d_2]$ with a note that at least one among $\alpha_{1,j}^{(u)}, \beta_{1,j}, \alpha_{2,j}^{(uv)}, \beta_{2,j}, \phi_{1,j}, \phi_{2,j}, \gamma_{1,j}^{(u)}, \gamma_{2,j}, \gamma_{3,j}$ and ξ_j is non-zero.

By applying the Fatou's lemma, we have

$$0 = \lim_{n \rightarrow \infty} \frac{\mathbb{E}_X[V(p_{G_n}(\cdot|X), p_{G_*}(\cdot|X))]}{m_n \mathcal{L}_2(G_n, G_*)} = \frac{1}{2} \int \liminf_{n \rightarrow \infty} \frac{|p_{G_n}(Y|X) - p_{G_*}(Y|X)|}{m_n \mathcal{L}_2(G_n, G_*)} d(X, Y),$$

which implies that $[p_{G_n}(Y|X) - p_{G_*}(Y|X)]/[m_n \mathcal{L}_2(G_n, G_*)] \rightarrow 0$ as $n \rightarrow \infty$ for almost surely (X, Y) . Since the term $\sum_{j=1}^{N^*} \exp(\log(1 + \exp((a_j^*)^\top X + b_j^*)))$ is bounded, we also have $T_n(Y|X)/[m_n \mathcal{L}_2(G_n, G_*)] \rightarrow 0$ as $n \rightarrow \infty$. Then, it follows that

$$0 = \lim_{n \rightarrow \infty} \frac{T_{n,1,1}(Y|X) + T_{n,1,2}(Y|X)}{m_n \mathcal{L}_2(G_n, G_*)} - \lim_{n \rightarrow \infty} \frac{T_{n,2,1}(Y|X) + T_{n,2,2}(Y|X)}{m_n \mathcal{L}_2(G_n, G_*)} + \lim_{n \rightarrow \infty} \frac{T_{n,3}(Y|X)}{m_n \mathcal{L}_2(G_n, G_*)}, \quad (23)$$

for almost surely $(X, Y) \in \mathcal{X} \times \mathcal{Y}$, where we have

$$\begin{aligned}
\lim_{n \rightarrow \infty} \frac{T_{n,1,1}(Y|X)}{m_n \mathcal{L}_2(G_n, G_*)} &:= \sum_{j \in [N^*]; |\mathcal{C}_j|=1} \left[\frac{\sum_{u=1}^d \alpha_{1,j}^{(u)} X^{(u)} + \phi_{1,j}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} F_{0,j}(Y|X) \right. \\
&\quad \left. + \left(\sum_{u=1}^d \alpha_{1,j}^{(u)} X^{(u)} + \phi_{1,j} \right) F_{1,j}(Y|X) + \frac{1}{2} \beta_{1,j} F_{2,j}(Y|X) \right], \\
\lim_{n \rightarrow \infty} \frac{T_{n,1,2}(Y|X)}{m_n \mathcal{L}_2(G_n, G_*)} &:= \sum_{j \in [N^*]; |\mathcal{C}_j|>1} \left[\left(\frac{\sum_{u=1}^d \alpha_{1,j}^{(u)} X^{(u)} + \phi_{1,j}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right. \right. \\
&\quad \left. \left. + \frac{\sum_{u,v=1}^d \frac{\alpha_{2,j}^{(uv)}}{1 + \mathbf{1}_{\{u=v\}}} X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} + \frac{\sum_{u=1}^d \gamma_{2,j}^{(u)} X^{(u)} + \frac{1}{2} \phi_{2,j}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right) F_{0,j}(Y|X) \right. \\
&\quad \left. + \left(\sum_{u=1}^d \alpha_{1,j}^{(u)} X^{(u)} + \phi_{1,j} + \frac{2 \sum_{u,v=1}^d \frac{\alpha_{2,j}^{(uv)}}{1 + \mathbf{1}_{\{u=v\}}} X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right. \right. \\
&\quad \left. \left. + \frac{\phi_{2,j} + 2 \sum_{u=1}^d \gamma_{2,j}^{(u)} X^{(u)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right) F_{1,j}(Y|X) + \left(\frac{1}{2} \beta_{1,j} + \sum_{u,v=1}^d \frac{\alpha_{2,j}^{(uv)} X^{(u)} X^{(v)}}{1 + \mathbf{1}_{\{u=v\}}} + \frac{1}{2} \phi_{2,j} \right. \right. \\
&\quad \left. \left. + \sum_{u=1}^d \gamma_{2,j}^{(u)} X^{(u)} + \frac{1}{2} \cdot \frac{\sum_{u=1}^d \gamma_{1,j}^{(u)} X^{(u)} + \gamma_{3,j}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right) F_{2,j}(Y|X) \right. \\
&\quad \left. + \left(\sum_{u=1}^d \frac{1}{2} \gamma_{1,j}^{(u)} X^{(u)} + \frac{1}{2} \gamma_{3,j} \right) F_{3,j}(Y|X) + \frac{1}{4} \beta_{2,j} F_{4,j}(Y|X) \right],
\end{aligned}$$

2322 and
2323

$$2324 \lim_{n \rightarrow \infty} \frac{T_{n,2,1}(Y|X)}{m_n \mathcal{L}_2(G_n, G_*)} := \sum_{j \in [N^*]: |\mathcal{C}_j|=1} \frac{\sum_{u=1}^d \alpha_{1,j}^{(u)} X^{(u)} + \phi_{1,j}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} H_j(Y|X),$$

2326

$$2327 \lim_{n \rightarrow \infty} \frac{T_{n,2,2}(Y|X)}{m_n \mathcal{L}_2(G_n, G_*)} := \sum_{j \in [N^*]: |\mathcal{C}_j|>1} \left[\frac{\sum_{u=1}^d \alpha_{1,j}^{(u)} X^{(u)} + \phi_{1,j}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right. \\ 2328 \\ 2329 \left. + \frac{\sum_{u,v=1}^d \frac{\alpha_{2,j}^{(uv)}}{1+1_{\{u=v\}}} X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} + \frac{\sum_{u=1}^d \gamma_{2,j}^{(u)} X^{(u)} + \frac{1}{2} \phi_{2,j}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \right] H_j(Y|X),$$

2330

2331

2332

2333

and

2334

2335

2336

$$\lim_{n \rightarrow \infty} \frac{T_{n,3}(Y|X)}{m_n \mathcal{L}_2(G_n, G_*)} := \sum_{j=1}^{N^*} \xi_j [F_{0,j}(Y|X) - H_j(Y|X)].$$

2337

It is worth noting that for almost every X , the set

2338

2339

2340

2341

$$\left\{ F_{\rho,j}(Y|X), H_j(Y|X) : 0 \leq \rho \leq 4, j \in [N^*] \right\}$$

2342

is linearly independent w.r.t Y . Therefore, it follows that the coefficients of those terms in the limit in equation (23) become zero.

2343

2344

For $j \in [N^*]$ such that $|\mathcal{C}_j| = 1$, by considering the coefficients of

2345

2346

2347

2348

2349

2350

2351

2352

- $F_{1,j}(Y|X)$, we have $\sum_{u=1}^d \alpha_{1,j}^{(u)} X^{(u)} + \phi_{1,j} = 0$ for almost surely X , indicating that $\alpha_{1,j}^{(u)} = \phi_{1,j} = 0$ for all $u \in [d]$;
- $F_{0,j}(Y|X)$, we have $\xi_j + \sum_{u=1}^d \alpha_{1,j}^{(u)} \cdot \frac{X^{(u)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} + \frac{\phi_{1,j}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} = 0$ for almost surely X . Since $\alpha_{1,j}^{(u)} = \phi_{1,j} = 0$ for all $u \in [d]$, we also get $\xi_j = 0$.
- $F_{2,j}(Y|X)$, we have $\beta_{1,j} = 0$.

2353

For $j \in [N^*]$ such that $|\mathcal{C}_j| > 1$, by considering the coefficients of

2354

2355

2356

2357

2358

2359

- $F_{1,j}(Y|X)$, we have
$$\sum_{u=1}^d \alpha_{1,j}^{(u)} X^{(u)} + \phi_{1,j} + \frac{2 \sum_{u,v=1}^d \frac{\alpha_{2,j}^{(uv)}}{1+1_{\{u=v\}}} X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} + \frac{\phi_{2,j} + 2 \sum_{u=1}^d \gamma_{2,j}^{(u)} X^{(u)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} = 0,$$

2360

for almost surely X . Since the set

2361

2362

2363

2364

2365

2366

$$\left\{ 1, X^{(u)}, \frac{1}{1 + \exp(-(a_j^*)^\top X - b_j^*)}, \frac{X^{(u)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)}, \right. \\ \left. \frac{X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} : u, v \in [d] \right\}$$

2367

2368

2369

2370

2371

2372

2373

2374

2375

- is linearly independent w.r.t X , we deduce $\alpha_{1,j}^{(u)} = \phi_{1,j} = \alpha_{2,j}^{(uv)} = \phi_{2,j} = \gamma_{2,j}^{(u)} = 0$ for all $u, v \in [d]$.
- $F_{0,j}(Y|X)$, we have
$$\xi_j + \frac{\sum_{u=1}^d \alpha_{1,j}^{(u)} X^{(u)} + \phi_{1,j}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} \\ + \frac{\sum_{u,v=1}^d \frac{\alpha_{2,j}^{(uv)}}{1+1_{\{u=v\}}} X^{(u)} X^{(v)}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} + \frac{\sum_{u=1}^d \gamma_{2,j}^{(u)} X^{(u)} + \frac{1}{2} \phi_{2,j}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} = 0,$$

- 2376 for almost surely X . Since $\alpha_{1,j}^{(u)} = \phi_{1,j} = \alpha_{2,j}^{(uv)} = \phi_{2,j} = \gamma_{2,j}^{(u)} = 0$ for all $u, v \in [d]$, we
 2377 get $\xi_j = 0$.
 2378
 2379 • $F_{3,j}(Y|X)$, we have $\sum_{u=1}^d \frac{1}{2} \gamma_{1,j}^{(u)} X^{(u)} + \frac{1}{2} \gamma_{3,j} = 0$ for almost surely X , indicating that
 2380 $\gamma_{1,j}^{(u)} = \gamma_{3,j} = 0$ for all $u \in [d]$;
 2381
 2382 • $F_{2,j}(Y|X)$, we have
 2383
$$\frac{1}{2} \beta_{1,j} + \sum_{u,v=1}^d \frac{\alpha_{2,j}^{(uv)} X^{(u)} X^{(v)}}{1 + \mathbb{1}_{\{u=v\}}} + \frac{1}{2} \phi_{2,j} + \sum_{u=1}^d \gamma_{2,j}^{(u)} X^{(u)} + \frac{1}{2} \frac{\sum_{u=1}^d \gamma_{1,j}^{(u)} X^{(u)} + \gamma_{3,j}}{1 + \exp(-(a_j^*)^\top X - b_j^*)} = 0,$$

 2384
 2385
 2386 for almost surely X . Since $\alpha_{2,j}^{(uv)} = \phi_{2,j} = \gamma_{2,j}^{(u)} = \gamma_{1,j}^{(u)} = \gamma_{3,j} = 0$ for all $u, v \in [d]$, we
 2387 also get $\beta_{1,j} = 0$.
 2388
 2389 • $F_{4,j}(Y|X)$, we have $\beta_{2,j} = 0$.
 2390

2391 Putting the above results together, we have $\xi_j = \alpha_{1,j}^{(u)} = \phi_{1,j} = \beta_{1,j} = \alpha_{2,j}^{(uv)} = \phi_{2,j} = \beta_{2,j} =$
 2392 $\gamma_{1,j}^{(u)} = \gamma_{2,j}^{(u)} = \gamma_{3,j} = 0$ for all $j \in [N^*]$ and $u, v \in [d]$. This contradicts the fact that at least one
 2393 among them is different from zero. Consequently, we achieve the local part in equation (21).
 2394

2395 K.3 PROOF OF PROPOSITION 3.1

2397 In this proof, we first present some fundamental results on the density estimation problem for
 2398 M-estimators in van de Geer (2000) in Appendix K.3.1, and then provide the main proof in Ap-
 2399 pendix K.3.2.

2401 K.3.1 PRELIMINARIES

2402 To streamline our discussion, let us introduce some necessary concepts from the empirical process
 2403 theory. In particular, let $\mathcal{P}_k(\Theta)$ be the set of all conditional densities with respect to mixing measures
 2404 in $\mathcal{G}_N(\Theta)$, i.e.
 2405

$$2406 \mathcal{P}_N(\Theta) := \{p_G(Y|X) : G \in \mathcal{G}_N(\Theta)\}.$$

2407 Additionally, we also consider two following variants of the set $\mathcal{P}_N(\Theta)$:
 2408

$$2409 \bar{\mathcal{P}}_k(\Theta) := \{p_{(G+G_*)/2}(Y|X) : G \in \mathcal{G}_N(\Theta)\},$$

$$2410 \bar{\mathcal{P}}_N^{1/2}(\Theta) := \{p_{(G+G_*)/2}^{1/2}(Y|X) : G \in \mathcal{G}_N(\Theta)\}.$$

2412 Next, we define for each $\delta > 0$ a Hellinger ball centered around the true conditional density $p_{G_*}(Y|X)$
 2413 and intersect with the set $\bar{\mathcal{P}}_N^{1/2}(\Theta)$ as below
 2414

$$2415 \bar{\mathcal{P}}_N^{1/2}(\Theta, \delta) := \{p^{1/2}(Y|X) \in \bar{\mathcal{P}}_N^{1/2}(\Theta) : h(p_G, p_{G_*}) \leq \delta\}.$$

2417 Moreover, the size of this Hellinger ball is quantified by the following term:
 2418

$$2419 \mathcal{J}_B(\delta, \bar{\mathcal{P}}_N^{1/2}(\Theta, \delta)) := \int_{\delta^2/2^{13}}^{\delta} H_B^{1/2}(t, \bar{\mathcal{P}}_N^{1/2}(\Theta, t), \|\cdot\|_2) dt \vee \delta, \quad (24)$$

2422 where $H_B(t, \bar{\mathcal{P}}_N^{1/2}(\Theta, t), \|\cdot\|_2)$ stands for the bracketing entropy of $\bar{\mathcal{P}}_N^{1/2}(\Theta, t)$ under the L^2 -norm,
 2423 and $t \vee \delta := \max\{t, \delta\}$. Now, we are ready to recall the results in van de Geer (2000).

2424 **Lemma K.1** (Theorem 7.4, van de Geer (2000)). *Take $\Psi(\delta) \geq \mathcal{J}_B(\delta, \bar{\mathcal{P}}_N^{1/2}(\Theta, \delta))$ such that $\Psi(\delta)/\delta^2$
 2425 is a non-increasing function of δ . Then, for a universal constant c and $\sqrt{n}\delta_n^2 \geq c\Psi(\delta_n)$, we achieve
 2426 that*

$$2427 \mathbb{P}\left(\mathbb{E}_X[h(p_{\hat{G}_n}(\cdot|X), p_{G_*}(\cdot|X))] > \delta\right) \leq c \exp(-n\delta^2/c^2),$$

2428 for any $\delta \geq \delta_n$.
 2429

Proof of Lemma K.1 is available in van de Geer (2000). Apart from this result, we also need to introduce the upper bounds of the covering number $N(\varepsilon, \mathcal{P}_N(\Theta), \|\cdot\|_\infty)$ and the bracketing entropy $H_B(\varepsilon, \mathcal{P}_N(\Theta), \|\cdot\|_2)$ as follows:

Lemma K.2. *Suppose that Θ is a bounded set, then we have for any $\varepsilon \in (0, 1/2)$ that*

$$(a) \log N(\varepsilon, \mathcal{P}_N(\Theta), \|\cdot\|_\infty) \lesssim \log(1/\varepsilon);$$

$$(b) H_B(\varepsilon, \mathcal{P}_N(\Theta), \|\cdot\|_2) \lesssim \log(1/\varepsilon).$$

Proof of Lemma K.2. Part (a). Recall that Θ is a compact set, then there exists an ε -cover, which we denote as $\bar{\Theta}_\varepsilon$. Moreover, it can be verified that $|\bar{\Theta}_\varepsilon| \leq \mathcal{O}(\varepsilon^{-(d_2+1)N})$. Next, for each mixing measure $G = \sum_{i=1}^N \delta_{(W_{e_i}, \nu_i)} \in \mathcal{G}_N(\Theta)$, we consider another one $\bar{G} = \sum_{i=1}^N \delta_{(\bar{W}_{e_i}, \bar{\nu}_i)}$, where $(\bar{W}_{e_i}, \bar{\nu}_i) \in \bar{\Theta}_\varepsilon$ is the closest point to (W_{e_i}, ν_i) in this set for any $i \in [N]$. Subsequently, we demonstrate that the set

$$\mathcal{Q} := \left\{ p_{\bar{G}}(Y|X) : (\bar{W}_{e_i}, \bar{\nu}_i) \in \bar{\Theta}_\varepsilon, \forall i \in [N] \right\}$$

is an ε -cover of the metric space $(\mathcal{P}_N(\Theta), \|\cdot\|_\infty)$. In other words, we need to show that for any $p_G(Y|X) \in \mathcal{P}_N(\Theta)$, there exists some density $p_{\bar{G}}(Y|X) \in \mathcal{Q}$ such that $\|p_G - p_{\bar{G}}\|_\infty \lesssim \varepsilon$.

Next, we decompose the term $T_n(Y|X) := \left[\sum_{j=1}^N \exp(\log(1 + \exp(g(X, \bar{W}_{e_j}))) \right] \cdot [p_G(Y|X) - p_{\bar{G}}(Y|X)]$ as

$$\begin{aligned} T_n(Y|X) &= \sum_{i=1}^N \exp(\log(1 + \exp(g(X, W_{e_i})))) \left[f(Y|g(X, W_{e_i}), \nu_i) - f(Y|g(X, \bar{W}_{e_i}), \bar{\nu}_i) \right] \\ &+ \sum_{i=1}^N \left[\exp(\log(1 + \exp(g(X, W_{e_i})))) - \exp(\log(1 + \exp(g(X, \bar{W}_{e_j})))) \right] \cdot \left[f(Y|g(X, \bar{W}_{e_i}), \bar{\nu}_i) - p_G(Y|X) \right]. \end{aligned}$$

As Θ and \mathcal{X} are bounded, we may assume that $\exp(\log(1 + \exp(g(X, W_{e_i})))) \leq B_1$ and $|f(Y|g(X, \bar{W}_{e_i}), \bar{\nu}_i) - p_G(Y|X)| \leq B_2$ for some positive constants B_1, B_2 . Thus, we obtain that

$$|T_n(Y|X)| \lesssim \sum_{i=1}^N B_1 \cdot \left[\|W_{e_i} - \bar{W}_{e_i}\| + |\nu_i - \bar{\nu}_i| \right] + \sum_{i=1}^N B_2 \cdot \|W_{e_i} - \bar{W}_{e_i}\| \lesssim \varepsilon.$$

Additionally, since the term $\sum_{j=1}^K \exp(|g(X, \bar{W}_{e_j})|)$ is bounded, we obtain $|p_G(Y|X) - p_{\bar{G}}(Y|X)| \lesssim \varepsilon$ for almost surely (X, Y) , or equivalently,

$$\|p_G - p_{\bar{G}}\|_\infty = \sup_{(X, Y) \in \mathcal{X} \times \mathcal{Y}} |p_G(Y|X) - p_{\bar{G}}(Y|X)| \lesssim \varepsilon.$$

This result indicates that \mathcal{Q} is an ε -cover of the metric space $(\mathcal{P}_N(\Theta), \|\cdot\|_\infty)$. Therefore, we get

$$N(\varepsilon, \mathcal{P}_N(\Theta), \|\cdot\|_\infty) \leq |\bar{\Theta}_\varepsilon| \leq \mathcal{O}(\varepsilon^{-(d_2+1)N}),$$

or equivalently,

$$\log N(\varepsilon, \mathcal{P}_N(\Theta), \|\cdot\|_\infty) \leq |\bar{\Theta}_\varepsilon| \lesssim \log(1/\varepsilon).$$

Part (b). Firstly, we will derive an upper bound for the Gaussian experts $f(Y|g(X, W_e), \nu)$. Since Θ is a compact set, we have $|g(X, W_e)| \leq M_1$ and $M_2 \leq \nu \leq M_3$ for any $X \in \mathcal{X}$ and $(W_e, \nu) \in \Theta$. Then, it follows that $f(Y|g(X, W_e), \nu) \leq B(Y|X)$, where

$$B(Y|X) := \begin{cases} \frac{1}{\sqrt{2\pi}M_2} \exp(-Y^2/(8M_3^2)), & \text{for } |Y| \geq 2M_1 \\ 1, & \text{for } |Y| < 2M_1, \end{cases}$$

for any $X \in \mathcal{X}$. Next, let $\eta \leq \varepsilon$ be some positive constant that we choose later, then we denote $\{\pi_1, \pi_2, \dots, \pi_N\}$ as an η -cover over $\mathcal{P}_N(\Theta)$. Based on this cover, we build the following brackets

2484 $L_i(Y|X) := \max\{\pi_i(Y|X) - \eta, 0\}$ and $U_i(Y|X) := \max\{\pi_i(Y|X) + \eta, B(Y|X)\}$, for any
 2485 $i \in [N]$. We can validate that $\mathcal{P}_N(Y|X) \subseteq \bigcup_{i=1}^N [L_i(Y|X), U_i(Y|X)]$ and $U_i(X, Y) - L_i(X, Y) \leq$
 2486 $\min\{2\eta, B(Y|X)\}$. As a result, we have

$$2487 \quad \|U_i - L_i\|_2 = \left(\int [U_i(Y|X) - L_i(Y|X)]^2 d(X, Y) \right)^{1/2} \leq 2\eta.$$

2488 The above result implies that

$$2489 \quad H_B(2\eta, \mathcal{P}_N(\Theta), \|\cdot\|_2) \leq \log N(\eta, \mathcal{P}_N(\Theta), \|\cdot\|_\infty) \lesssim \log(1/\eta).$$

2490 Then, by setting $\eta = \varepsilon/2$, we arrive at

$$2491 \quad H_B(\varepsilon, \mathcal{P}_N(\Theta), \|\cdot\|_1) \lesssim \log(1/\varepsilon).$$

2492 Hence, the proof is completed. \square

2493 K.3.2 MAIN PROOF

2494 Since $\bar{\mathcal{P}}_N^{1/2}(\Theta, t) \subset \bar{\mathcal{P}}_N^{1/2}(\Theta)$ for any $t > 0$, we have

$$2495 \quad H_B(t, \bar{\mathcal{P}}_N^{1/2}(\Theta, t), \|\cdot\|_2) \leq H_B(t, \bar{\mathcal{P}}_N^{1/2}(\Theta), \|\cdot\|_2) = H_B(t/\sqrt{2}, \bar{\mathcal{P}}_N(\Theta), h), \quad (25)$$

2496 where the last equality is due to the relationship between the Hellinger distance h and the L^2 -norm.
 2497 Note that for any two mixing measure G and G' , Lemma 4.2 in van de Geer (2000) indicates that

$$2498 \quad h^2\left(\frac{1}{2}p_G + \frac{1}{2}p_{G_*}, \frac{1}{2}p_{G'} + \frac{1}{2}p_{G_*}\right) \leq \frac{1}{2}h^2(p_G, p_{G'}),$$

2499 which yields $H_B(t/\sqrt{2}, \bar{\mathcal{P}}_N(\Theta), h) \leq H_B(t, \mathcal{F}_{k_1, k_2}(\Theta), h)$. This result together with equation equa-
 2500 tion (25) implies that

$$2501 \quad H_B(t, \bar{\mathcal{P}}_N^{1/2}(\Theta, t), \|\cdot\|_2) \leq H_B(t, \mathcal{P}_N(\Theta), h).$$

2502 From equation (24) and part (b) of Lemma K.2, we have that

$$2503 \quad \begin{aligned} \mathcal{J}_B(\delta, \bar{\mathcal{P}}_N^{1/2}(\Theta, \delta)) &= \int_{\delta^2/2^{13}}^{\delta} H_B^{1/2}(t, \bar{\mathcal{P}}_N^{1/2}(\Theta, t), \|\cdot\|_2) dt \vee \delta \\ &\leq \int_{\delta^2/2^{13}}^{\delta} H_B^{1/2}(t, \bar{\mathcal{P}}_N^{1/2}(\Theta, t), h) dt \vee \delta \\ &\lesssim \int_{\delta^2/2^{13}}^{\delta} \log(1/t) dt \vee \delta. \end{aligned}$$

2504 Next, let $\Psi(\delta) = \delta\sqrt{\log(1/\delta)}$, then it can be verified that $\Psi(\delta)/\delta^2$ is a non-increasing function of δ .

2505 Furthermore, the above result indicates that $\Psi(\delta) \geq \mathcal{J}_B(\delta, \bar{\mathcal{F}}_{k_1, k_2}^{1/2}(\Theta, \delta), \|\cdot\|_2)$. By considering the
 2506 sequence (δ_n) defined as $\delta_n := \sqrt{\log(n)/n}$, we have $\sqrt{n}\delta_n^2 \geq c\Psi(\delta_n)$ for some universal constant
 2507 $c > 0$. It follows from Lemma K.1 that

$$2508 \quad \mathbb{P}\left(\mathbb{E}_X[h(p_{\hat{G}_n}(\cdot|X), p_{G_*}(\cdot|X))] > C\sqrt{\log(n)/n}\right) \lesssim \exp(-c\log(n)),$$

2509 for some universal constant $C > 0$ depending only on Θ . Since the Total Variation distance is upper
 2510 bounded by the Hellinger distance, we deduce

$$2511 \quad \mathbb{P}\left(\mathbb{E}_X[V(p_{\hat{G}_n}(\cdot|X), p_{G_*}(\cdot|X))] > C\sqrt{\log(n)/n}\right) \lesssim \exp(-c\log(n)),$$

2512 or equivalently,

$$2513 \quad \mathbb{E}_X[V(p_{\hat{G}_n}(\cdot|X), p_{G_*}(\cdot|X))] = \mathcal{O}_P(\sqrt{\log(n)/n}).$$

2514 Hence, the proof is completed.

2538 L LLM USAGE

2539

2540 We affirm that LLMs did not play a significant role in the development of this work, to the extent that
2541 they could be regarded as an author. No content generation, ideation, or technical writing assistance
2542 was delegated to LLMs.

2543

2544 M BROADER IMPACT

2545

2546 Although our work mostly contributes to the machine learning literature, it also draws inspiration from
2547 biology and neuroscience. Specifically, the competition mechanism is rooted in biology, has been
2548 studied in neuroscience, and has motivated a few machine learning algorithms. Our work contributed
2549 a theoretically grounded algorithm to train large-scale SMOE models, which could potentially push
2550 the frontier of the next LLM generation. Lastly, working with large models requires rather costly
2551 resources. We took serious precautions during the development of this work, including providing a
2552 guideline for hyper-parameter selection, and conducting a single experiment using the same random
2553 seed to ensure the results are reliable at a low cost.

2554

2555

2556

2557

2558

2559

2560

2561

2562

2563

2564

2565

2566

2567

2568

2569

2570

2571

2572

2573

2574

2575

2576

2577

2578

2579

2580

2581

2582

2583

2584

2585

2586

2587

2588

2589

2590

2591