

# MIXTURE OF LoRA EXPERTS

Xun Wu<sup>1,2,\*</sup>, Shaohan Huang<sup>1,✉</sup>, Furu Wei<sup>1</sup>

<sup>1</sup>Microsoft Research Asia <sup>2</sup>Tsinghua University  
 wuxun21@mails.tsinghua.edu.cn; {shaohanh, fuwei}@microsoft.com

## ABSTRACT

Low-Rank Adaptation (LoRA) (Hu et al., 2021) has emerged as a pivotal technique for fine-tuning large pre-trained models, renowned for its efficacy across a wide array of tasks. The modular architecture of LoRA has catalyzed further research into the synergistic composition of multiple trained LoRAs, aiming to amplify performance across various tasks. However, the effective composition of these trained LoRAs presents a formidable challenge: (1) Linear arithmetic composition can lead to the diminution of the generative capabilities inherent in the original pre-trained models or the distinctive attributes of the individually trained LoRAs, potentially resulting in suboptimal outcomes. (2) Reference tuning-based composition exhibits limitations in adaptability and incurs significant computational costs due to the requirements to retrain a large model. In response to these challenges, we propose **Mixture of LoRA Experts (MOLE)**. MOLE treats each layer of trained LoRAs as a distinct expert and implements hierarchical weight control by integrating a learnable gating function within each layer to learn optimal composition weights tailored specifically to the objectives of a given domain. MOLE not only demonstrates enhanced performance in LoRA composition but also preserves the essential flexibility necessary for effective composition of trained LoRAs with minimal computational overhead. Extensive experiments conducted in both Natural Language Processing (NLP) and Vision & Language (V&L) domains validate the effects of MOLE. Our code are available at <https://github.com/yushuiwx/MOLE.git>.

## 1 INTRODUCTION

Recent advances in deep learning have been driven by large-scale pre-trained models such as OPT (Zhang et al., 2022), LLaMA (Touvron et al., 2023) in the Natural Language Processing (NLP) domain and CLIP (Radford et al., 2021a), DALL·E 2 (Ramesh et al., 2022) in the Vision & Language (V&L) domain. These models show outstanding performance across various tasks when fine-tuned on down-stream datasets, but their increasing size entails significant computational costs for full fine-tuning. To mitigate this, LoRA (Hu et al., 2021) is introduced. By freezing the pretrained model weights and injecting trainable rank decomposition matrices, LoRA is proven to be an effective fine-tuning methodology in scenarios with constrained computational resources (Lester et al., 2021; An et al., 2022).

While LoRA serves as plug-and-play plugins for pre-trained models, recent initiatives explores the composition of separate trained LoRAs to achieve joint generation of learned characteristics (Huang et al., 2023; Zhang et al., 2023; Ruiz et al., 2023). However, these efforts may encounter several challenges. As shown in Figure 2 (a), linear arithmetic composition (Zhang et al., 2023; Huang et al., 2023; Han et al., 2023) composes trained LoRAs

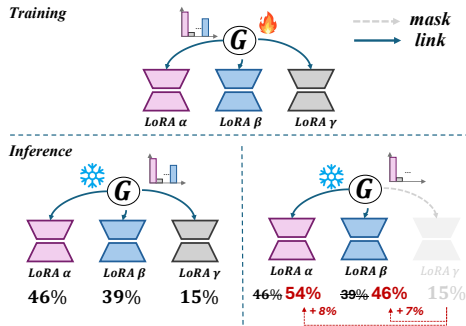


Figure 1: **Workflow of MOLE**. In the training phase, MOLE predicts weights for multiple LoRAs. In the inference phase, MOLE can allocate weights to multiple LoRAs, or, without altering the gating weights, achieve a more flexible LoRA composition by masking out undesired LoRAs and recalculating and distributing weights proportionally.

\*Contribution during internship at Microsoft. ✉ Corresponding Author.

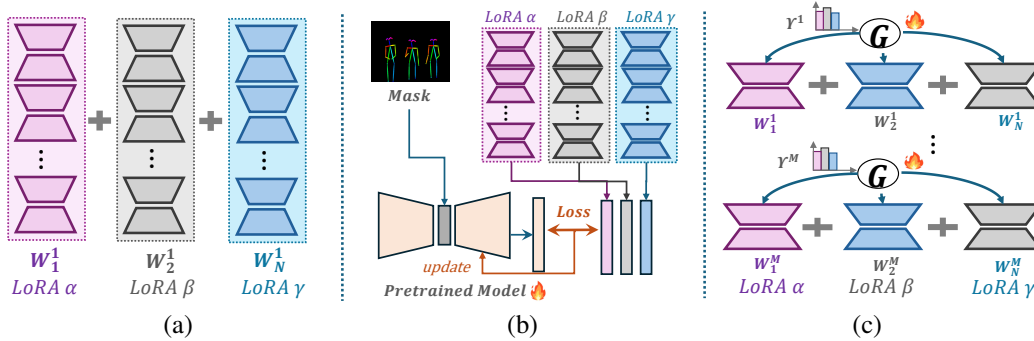


Figure 2: Overview of LoRA composition methods: (a) Linear arithmetic composition (Eq.2), which commonly applies the same composition weight  $W_i$  to all layers of the  $i^{th}$  LoRA. (b) Reference tuning-based composition involves retraining a large model by integrating outputs from multiple LoRAs using manually-crafted mask information. (c) Our MOLE, which learns a distribution  $\Upsilon^j$  for the  $j^{th}$  layer of LoRAs to determine the composition weight  $W_i^j$ .

directly. However, composing multiple LoRAs (typically  $\geq 3$ ) can impair the generative performance of pre-trained models. To mitigate this, weight normalization is applied prior to the composition, but may erase the unique characteristics of individual trained LoRAs as the composing weight of each LoRA is reduced (refer to Observation 1 in § 3.1). Another approach, as depicted in Figure 2 (b), known as reference tuning-based composition (Gu et al., 2023), is tailored for the V&L domain and achieves superior performance. However, it is limited in terms of LoRA flexibility due to the utilization of manually-designed masks and involves substantial training costs, necessitating a full model retraining. In light of the above situation, an important question arises:

*How can multiple trained LoRAs be composed dynamically and efficiently, while preserving all their individual characteristics?*

To address that issues, we introduce **Mixture of LoRA Experts (MOLE)**. Recognizing that individual layers of a trained LoRA exhibit distinct characteristics, which collectively define the overall characteristic of the trained LoRA (refer to Observation 2 in § 3.1), MOLE involves modulating the weights of different trained LoRAs within each layer, which we refer to as “hierarchical weight contro”. As shown in Figure 2 (c), MOLE views each layer of trained LoRAs as a individual expert and incorporates a gating function within each layer to learn the optimal composition weights based on a specified domain objective. This dynamically enhances desirable characteristics while mitigating less favorable ones, ultimately achieving a more effective composition of LoRAs and prevents the loss of desirable LoRA characteristics that may occur in linear arithmetic composition.

Additionally, unlike reference tuning-based composition (Gu et al., 2023), our MOLE maintains flexibility in composing multiple trained LoRAs with reduced computational costs. As the workflow of MOLE shown in Figure 1, during training, MOLE learns the gating function for multiple trained LoRAs and keep all other parameters frozen, resulting in minimal computational costs. During inference, MOLE has two inference modes: In the first mode, MOLE utilizes all trained LoRAs with the learned gating function, preserving their individual characteristics with allocated weights. During the second mode, MOLE allows manual masking of unwanted LoRAs and recalculates and distributes weights proportionally without the need for retraining. These two modes enable MOLE to adapt to different scenarios, providing a versatile and flexible approach for effective LoRA composition.

We validate the effects of MOLE in both NLP and V&L domains. Our findings, encompassing both qualitative and quantitative results, demonstrate that MOLE outperforms existing LoRA composition approaches. The contributions of our paper are the following:

- We introduce a significant and intricate problem: how to dynamically and efficiently compose multiple trained LoRAs while preserving all their individual characteristics, to further investigate the applicability of LoRA in real-world scenarios.

- We introduce Mixture of LoRA Experts (MOLE), a method that achieves a more efficient and flexible composition of multiple trained LoRAs by employing hierarchical weight control through learnable gating functions within each layer of trained LoRAs.
- Extensive experiments on both V&L and NLP domain demonstrate that MOLE can enhance LoRA composition performance and mitigates issues associated with existing composition methods.

## 2 BACKGROUND

### 2.1 LoRAs COMPOSITION

LoRA (Hu et al., 2021) is a parameter-efficient fine-tuning method to adapt large models to novel tasks and shows superior performance (Hu et al., 2021; Huang et al., 2023; Zhang et al., 2023; Sung et al., 2022). In practical applications, a individual LoRA often fall short of meeting user expectations. A common solution is to compose multiple trained LoRAs, each specialized in specific aspects (e.g., clothing or facial features), with the aim of creating a comprehensive character representation. Research on LoRA composition is limited and primarily concentrates on two distinct methodologies as follows:

**Linear arithmetic composition.** As shown in Figure 2 (a), the most commonly employed composition method is directly composing multiple LoRAs, i.e.,

$$\hat{\mathbf{W}} = \mathbf{W} + \sum_{i=1}^N \Delta \mathbf{W}_i, \quad (1)$$

where  $\mathbf{W}$  indicates the original parameter of pre-trained model and  $\Delta \mathbf{W}_i$  denotes the  $i^{th}$  trained LoRA. However, this manner may affect the original weight  $\mathbf{W}$  when  $N$  increasing, thereby diminishing the model’s generative capabilities. So, it is common practice to normalize the composition weights, termed as normalized linear arithmetic composition, i.e.,

$$\hat{\mathbf{W}} = \mathbf{W} + \sum_{i=1}^N w_i \cdot \Delta \mathbf{W}_i, \quad (2)$$

where  $\sum_{i=1}^N w_i = 1$ . This manner prevents any adverse impact on the embedding of the original model, but leading to the loss of individual LoRA characteristics, as the composing weight  $w_i$  for each trained LoRA is reduced (Gu et al., 2023).

In NLP domain, PEMs (Zhang et al., 2023) first define arithmetic operators for LoRA, and explore the effectiveness of composing multiple LoRAs in several scenarios. LoRAhub (Huang et al., 2023) utilizes a gradient-free manner to estimate the composition weights of trained LoRAs and achieves adaptable performance on unseen tasks. In V&L domain, SVDiff (Han et al., 2023) introduces a arithmetic-based manner to compose multiple visual concepts into a single image.

**Reference tuning-based composition.** As shown in Figure 2 (b), reference tuning-based composition (Gu et al., 2023) tackles the limitations of linear arithmetic composition by introducing gradient fusion and controllable sampling. However, it suffers from compositional inflexibility due to manually designed masks, which necessitates retraining when incorporating different LoRAs or creating new masks. Moreover, this approach entails retraining large models, resulting in substantial computational costs.

It is important to note that reference tuning-based composition relies on position masks, which distinguishes it from our model. Consequently, direct comparisons may not be appropriate due to the fundamentally different underlying principles. Therefore, our primary focus in this paper is to compare MOLE with linear arithmetic composition.

### 2.2 MIXTURE-OF-EXPERTS

Mixture-of-Experts (MoE) (Xie et al., 2023) is a promising approach to scale up the number of parameters within the same computational bounds. Different from standard transformer models, each MoE layer consists of  $N$  independent feed-forward networks  $\{\mathbf{E}_i\}_{i=0}^N$  as the experts, along with a

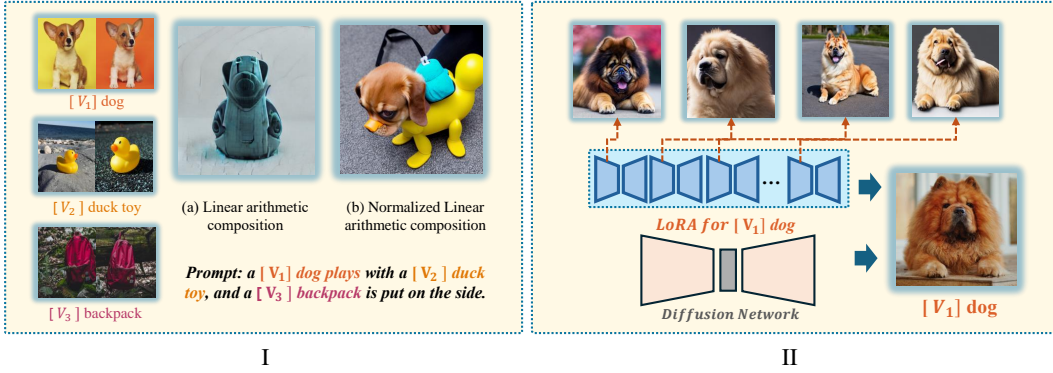


Figure 3: I. Results of (a) linear arithmetic composition (Eq. 1) and (b) normalized linear arithmetic composition (Eq. 2) based on Dreambooth (Ruiz et al., 2023). II. Visualization of the effects for different layers in LoRA by selectively activating specific parameters from the network, moving from the beginning to the end.

gating function  $\alpha(\cdot)$  to model a probability distribution indicating the weights over these experts' outputs. For the hidden representation  $\mathbf{h} \in \mathbb{R}^d$  of input token, the gate value of routing  $\mathbf{h}$  to expert  $\mathbf{E}_i$  is denoted as:

$$\alpha(\mathbf{E}_i) = \exp(\mathbf{h} \cdot \mathbf{e}_i) / \sum_{j=0}^N \exp(\mathbf{h} \cdot \mathbf{e}_j), \quad (3)$$

where  $\mathbf{e}_i$  denotes the trainable embedding of  $\mathbf{E}_i$ . Then, the corresponding  $k$  experts, according to the top- $k$  gated values, are activated and the output  $\mathbf{O}$  of the MoE layer is

$$\mathbf{O} = \mathbf{h} + \sum_{i=0}^N \alpha(\mathbf{E}_i) \cdot \mathbf{E}_i(\mathbf{h}). \quad (4)$$

### 3 METHOD

In this section, we first introduce some motivating observations in § 3.1. Then, we introduce the structure details and training objectives of MOLE in § 3.2 and § 3.3, respectively.

#### 3.1 MOTIVATING OBSERVATION

**Observation 1:** Directly composing multiple trained LoRAs (Eq. 1) impacts the model's generative ability, whereas applying weight normalization (Eq. 2) preserves this capacity but may sacrifice LoRA characteristics.

Specifically, in V&L domain, as depicted in Figure 3 I, we observe that directly composing multiple trained LoRAs into the original embedding led to significant parameter variations, resulting in meaningless output. Furthermore, when normalization was applied, some of the original characteristics of these trained LoRAs are indeed compromised. These observations align with those elaborated upon in (Gu et al., 2023).

In NLP domain, when composing four or more LoRAs within the FLAN-T5 (Chung et al., 2022) model, we observed that the model's output became disordered. Furthermore, implementing weight normalization for LoRAs trained across five datasets, as presented in Table 4, led to a decreased performance of the composition model. This suggests that while weight normalization preserves generative capacity, it adversely affects the intrinsic qualities of these trained LoRAs.

**Observation 2:** Individual layers of a trained LoRA exhibit unique traits, which cumulatively define the LoRA's overall attributes.

Inspired by the findings of (Voynov et al., 2023), which revealed that different layers in text-to-image models govern various attributes, such as style and color, we investigate the features learned

by different layers within LoRA. In V&L domain, as illustrated in Figure 3 II, we observed that different layers of LoRA encode distinct features, such as dog coat color and facial features. In NLP domain, we trained a single LoRA on a combined dataset comprising ANLI-R1 (Nie et al., 2019), ANLI-R2 (Nie et al., 2019), and QNLI (Rajpurkar et al., 2018) datasets, as depicted in Table 5. Notably, when evaluated on these sub-datasets, we observed significant variations in performance across different layers of this LoRA. Specifically, the layers ranging from 0% to 20% performed best on QNLI, the layers spanning from 40% to 60% excelled on ANLI-R2, and the layers covering 80% to 100% outperformed others on ANLI-R1. This observation inspires that we can dynamically optimize the layer-specific weights according to a defined domain objective, enhancing desirable characteristics while suppressing less favorable ones, thereby achieving a more effective composition of trained LoRAs.

### 3.2 MIXTURE OF LORA EXPERTS

Drawing inspiration from above observations, we introduce the Mixture of LoRA Experts.

Referring to Figure 4, consider a transformer block within the pre-trained model, parameterized by  $\theta$  (encompassing both the multi-head attention layer and the feed-forward neural network), and a set of corresponding trained LoRAs  $\Omega = \{\Delta\theta_i\}_{i=0}^N$  where  $N$  represents the number of trained LoRA candidates, when given an input  $\mathbf{x} \in \mathbb{R}^{L \times d}$ , the output of the pre-trained model block  $\theta$  is presented as  $\mathbf{F}_\theta \in \mathbb{R}^{L \times d}$ :

$$\mathbf{x}'_\theta = \mathbf{x} + f_{\text{Attn}}(\text{LN}(\mathbf{x})|\theta), \quad (5)$$

$$\mathbf{F}_\theta(\mathbf{x}) = \mathbf{x}'_\theta + f_{\text{FFN}}(\text{LN}(\mathbf{x}'_\theta)|\theta), \quad (6)$$

where  $L$  and  $d$  indicate the sequence length and the dimension of  $\mathbf{x}$ , respectively.  $f_{\text{Attn}}(\cdot)$  and  $f_{\text{FFN}}(\cdot)$  denotes the multi-head attention layer and feed-forward neural network, respectively. LN refers to layer normalization. The output of each LoRA is presented as  $\mathbf{E}_{\Delta\theta_i}(\mathbf{x}) \in \mathbb{R}^{L \times d}$ ,

$$\mathbf{x}'_{\Delta\theta_i} = \mathbf{x} + f_{\text{Attn}}(\text{LN}(\mathbf{x})|\Delta\theta_i), \quad (7)$$

$$\mathbf{E}_{\Delta\theta_i}(\mathbf{x}) = \mathbf{x}'_{\Delta\theta_i} + f_{\text{FFN}}(\text{LN}(\mathbf{x}'_{\Delta\theta_i})|\Delta\theta_i). \quad (8)$$

After that, MOLE applies a learnable gating function  $\mathcal{G}(\cdot)$  to model the optimal distribution of composition weights for outputs of these trained LoRAs. Specifically, by taking  $\{\mathbf{E}_{\Delta\theta_i}(\mathbf{x})\}_{i=0}^N$  as input,  $\mathcal{G}(\cdot)$  first apply concatenation (denoted as  $\oplus$ ) and normalization (for training stability), i.e.

$$\mathbf{E}_\Omega(\mathbf{x}) = \text{Normalization}(\mathbf{E}_{\Delta\theta_0}(\mathbf{x}) \oplus \dots \oplus \mathbf{E}_{\Delta\theta_{N-1}}(\mathbf{x})), \quad (9)$$

where  $\mathbf{E}_\Omega(\mathbf{x}) \in \mathbb{R}^\xi$  and  $\xi = N \times L \times d$ .  $\oplus$  indicates the concatenation operation. Then we flatten and reduce the  $\mathbf{E}_\Omega(\mathbf{x})$  to  $N$ -dimensions by a dot-product operation with the learnable parameter  $\mathbf{e} \in \mathbb{R}^{\xi \times N}$  in the gating function  $\mathcal{G}(\cdot)$ ,

$$\varepsilon = \text{Flatten}(\mathbf{E}_\Omega(\mathbf{x}))^\top \cdot \mathbf{e}, \quad \varepsilon \in \mathbb{R}^N, \quad (10)$$

The gate value for each LoRA is computed as

$$\mathcal{G}(\varepsilon_i) = \frac{\exp(\varepsilon_i/\tau)}{\sum_{j=1}^N \exp(\varepsilon_j/\tau)}, \quad (11)$$

the temperature scalar  $\tau$  is learnable. The final output  $\tilde{\mathbf{E}}_\Omega(\mathbf{x})$  of the gating function  $\mathcal{G}(\cdot)$  is obtained by multiplying the output of each LoRA expert with the corresponding gating values, presented as

$$\tilde{\mathbf{E}}_\Omega(\mathbf{x}) = \sum_{i=0}^N \mathcal{G}_i(\varepsilon_i) \cdot \mathbf{E}_{\Delta\theta_i}(\mathbf{x}), \quad (12)$$

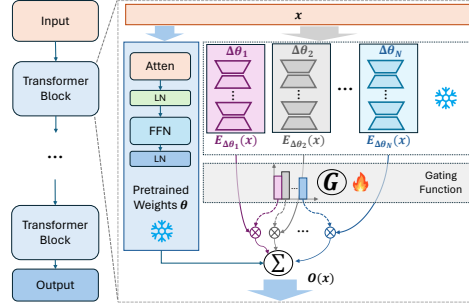


Figure 4: **Illustration of proposed MOLE.** MOLE employs a learnable gating function that utilizes the outputs of multiple LoRAs at each layer to determine composition weights.

Table 1: Text-alignment and image-alignment results for multiple LoRAs composition in CLIP feature space. NLA denotes normalized linear arithmetic composition (Eq. 2). The best performance is in bold.

# Visual Concepts	Text-alignment			Image-alignment, (Concept 1)			Image-alignment, (Concept 2)			Image-alignment, (Concept 3)		
	NLA	SVDiff	MoLE	NLA	SVDiff	MoLE	NLA	SVDiff	MoLE	NLA	SVDiff	MoLE
Fancy boot + Monster + Clock	0.754	0.742	0.832	0.781	0.758	0.784	0.791	0.749	0.801	0.763	0.812	0.809
Emoji + Car + Cartoon	0.610	0.607	0.696	0.619	0.734	0.839	0.711	0.702	0.709	0.652	0.686	0.679
Vase + Wolf plushie + Teapot	0.752	0.812	0.863	0.687	0.807	0.835	0.705	0.782	0.746	0.653	0.694	0.721
White Cat + Wolf plushie + Can	0.704	0.772	0.780	0.801	0.804	0.802	0.678	0.763	0.825	0.650	0.729	0.714
Shiny sneaker + Wolf plushie + Teapot	0.778	0.789	0.791	0.812	0.783	0.690	0.723	0.751	0.790	0.688	0.676	0.721
Car + Wolf plushie + Teapot	0.635	0.681	0.684	0.652	0.763	0.713	0.601	0.664	0.745	0.685	0.612	0.707
Can + Wolf plushie + backpack	0.601	0.782	0.754	0.653	0.705	0.767	0.602	0.755	0.782	0.681	0.738	0.723
Golden Retriever + Wolf plushie + Teapot	0.670	0.716	0.784	0.713	0.784	0.790	0.601	0.802	0.809	0.678	0.761	0.748
Golden Retriever + Boot + Monster	0.614	0.762	0.755	0.665	0.662	0.620	0.748	0.832	0.862	0.723	0.719	0.735
Backpack dog + Bowl + Teapot	0.607	0.712	0.703	0.653	0.672	0.756	0.734	0.720	0.755	0.692	0.688	0.701
Backpack dog + White Cat + Emoji	0.648	0.703	0.717	0.674	0.692	0.812	0.719	0.741	0.701	0.742	0.720	0.796
Dog + Wolf + Backpack	0.717	0.738	0.722	0.547	0.565	0.552	0.679	0.681	0.707	0.766	0.795	0.831
Cat + Sunglasses + Boot	0.770	0.791	0.837	0.845	0.793	0.815	0.845	0.793	0.815	0.845	0.793	0.815
Table + Can + Teapot	0.836	0.827	0.810	0.753	0.770	0.741	0.751	0.799	0.806	0.818	0.771	0.829
Robot + Dog + Clock	0.663	0.638	0.693	0.689	0.764	0.797	0.645	0.674	0.710	0.661	0.715	0.717
Average	0.678	0.728	<b>0.759</b>	0.715	0.746	<b>0.783</b>	0.682	0.731	<b>0.756</b>	0.686	0.708	<b>0.732</b>

in which  $\tilde{E}_\Omega(\mathbf{x}) \in \mathbb{R}^{L \times d}$  and  $\mathcal{G}_i(\cdot)$  represents the weight of the  $i^{th}$  trained LoRA. So, the final output of this block is computed by adding the output of the gating function to the output of the pre-trained network:

$$\mathbf{O}(\mathbf{x}) = \mathbf{F}_\theta(\mathbf{x}) + \tilde{E}_\Omega(\mathbf{x}). \quad (13)$$

Besides, we conducted an exploration of MOLE’s performance when employing gating functions at different hierarchical levels (layer-wise and matrix-wise, etc). Please refer to Section 5.

### 3.3 TRAINING OBJECTIVE

**Gating Balancing Loss.** As shown in Figure 5 (a), we observed that the average entropy of the distribution probabilities from the gating functions gradually decreases as the number of training steps increases, i.e., the gating function tends to converge to a state where it always produces large weights for a early-stage well-performing LoRA (e.g., shown in Figure. 5 (b), 68% gating probability for LoRA  $\beta$  among three LoRAs), leading to only a handful of LoRAs having a significant impact in the end and a loss of the characteristics of other LoRAs. To alleviate this, we propose a gating balancing loss  $\mathcal{L}_{\text{balance}}$  as

$$\mathcal{L}_{\text{balance}} = -\log \left( \prod_{i=0}^N \mathbf{q}^{(i)} \right), \quad (14)$$

where

$$\mathbf{q}^{(i)} = \frac{1}{M} \sum_{k=1}^M \frac{\exp(\varepsilon_i^k / \tau)}{\sum_{j=1}^N \exp(\varepsilon_j^k / \tau)}, \quad (15)$$

and  $M$  represents the number of blocks where gating functions are placed and  $N$  denotes the number of LoRAs. This balanced loss encourages balanced gating because it is minimized when the dispatching is ideally balanced.

**Domain-specific Loss.** Additionally, for adaptation to different domains, we employ distinct domain-specific training objectives denoted as  $\mathcal{L}_D$ . In V&L domain, we employ unsupervised training with both local and global guidance from CLIP (Radford et al., 2021b) to optimize MoLE. In NLP domain, we follow the loss function in FLAN-T5 (Chung et al., 2022). The overall training objective  $\mathcal{L}$  is the weighted sum of the above-mentioned two losses, represented as:

$$\mathcal{L} = \mathcal{L}_D + \alpha \mathcal{L}_{\text{balance}}, \quad (16)$$

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

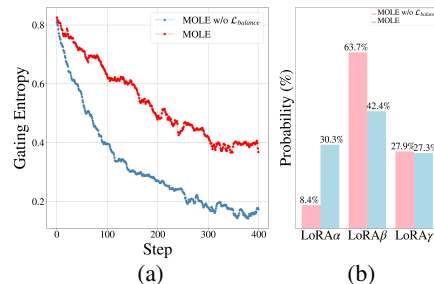


Figure 5: (a) The average gating entropy of all gating functions varies with the training steps. (b) The average weight distribution (%) of three LoRAs w and w/o  $\mathcal{L}_{\text{balance}}$ .

Table 2: Text-alignment and image-alignment results for multiple LoRA experts composition in CLIP feature space. The best performance is in bold and the second-best value is indicated with an underline. NLA denotes normalized linear arithmetic composition (Eq. 2). *SOTA full-parameter training methods are highlighted by*.

# Number of Concepts	Text-alignment					Average Image-alignment				
	NLA	Custom	Textual Inversion	SVDiff	MoLE	NLA	Custom	Textual Inversion	SVDiff	MoLE
3	0.678	<u>0.751</u>	0.709	0.728	<b>0.759</b>	0.694	<b>0.761</b>	0.720	0.719	<u>0.757</u>
4	0.681	<b>0.735</b>	0.721	0.717	<u>0.725</u>	0.712	<b>0.760</b>	0.736	0.721	<u>0.742</u>
5	0.652	<u>0.731</u>	0.704	0.723	<b>0.762</b>	0.682	<b>0.798</b>	0.710	0.708	<u>0.737</u>
6	0.678	0.722	<b>0.735</b>	0.709	<u>0.727</u>	0.698	0.721	<b>0.747</b>	0.712	<u>0.736</u>
Average	0.672	<u>0.734</u>	0.717	0.719	<b>0.752</b>	0.692	<b>0.760</b>	0.728	0.715	<u>0.743</u>

**Optimization Gating Function Only.** We freeze all trained LoRAs and pre-trained model parameters, optimizing only the gating function’s parameters. This helps preserve characteristics of trained LoRAs, particularly when training data is limited.

## 4 EXPERIMENTS

### 4.1 MOLE ON V&L DOMAIN

**Experimental Setup.** For V&L domain, we apply MOLE to multi-subjects text-to-image generation task and choose DreamBooth (Ruiz et al., 2023) (built on Stable Diffusion V2.1) as the base generator. Following the common setting (Han et al., 2023; Gal et al., 2022a), where 2 to 3 concepts are typically composed into a new multi-concept image, we conduct experiments by composing three separate trained LoRAs. During training MOLE, we process the image resolution to  $512 \times 512$  and set learning rate as  $1e-5$ . We use DDPM sampler (Ho et al., 2020) with 50 steps in each case and train 400 iterations for each required composition with batch size 2 and  $\alpha$  as 0.5.

**Metrics and Compared Baselines.** Following (Ruiz et al., 2023; Han et al., 2023), we evaluate our method on (1) Image alignment. The visual similarity of generated images with the individual composed concepts, using similarity in CLIP (Radford et al., 2021a) image feature space, (2) Text-alignment of the generated images with given text prompts, using text-image similarity in CLIP feature space (Radford et al., 2021a). For each composition, we calculated the average scores among 200 generated images per prompt using 5 text prompts. We compared our MOLE with normalized linear arithmetic composition (Eq. 2) and SVDiff (Han et al., 2023). Additionally, to further validate the effectiveness of MOLE, we also compare MOLE with state-of-the-art multi-subjects generation methods (full-parameters training based), which can be found in Section 5.

**Main Results.** As shown in Table 1, this study involves 15 different compositions of three visual subjects. The overall results show that our method significantly outperforms other comparative methods in terms of Text-alignment score, with a 0.031 average improvement compared to SVDiff, as well as the Image-alignment score associated with three visual concepts (e.g., 0.037 average improvement compared to SVDiff in concept 1), providing evidence of our MOLE’s superior capability in accurately capturing and depicting the subject information of user-provided images, as well as displaying multiple entities concurrently within a single image. Significantly, prior research (Kumari et al., 2023; Gal et al., 2022b) indicates a trade-off between Text-alignment and Image-alignment scores in multi-subjects generation. Excelling in both scores is challenging, highlighting the strength of our MOLE. Additionally, as shown in Figure 9, 10 and 11, our approach outperforms two other methods in preserving subject fidelity in generated images. The comparative methods often omit a subject, as seen in the NLA composition’s failure to include elements like “cat” in Figure 9 (line 2) and “barn” in Figure 10, and SVDiff’s inability to precisely represent “dog” and “cat” in Figure 10. Furthermore, while these methods can generate images with three subjects, there’s a noticeable leakage and mixing of appearance features, resulting in lower subject fidelity compared to user-provided images. In contrast, our method effectively retains the subjects specified by the user, with each accurately depicted.

### 4.2 MOLE ON NLP DOMAIN

**Experimental Setup.** For NLP domain, following (Huang et al., 2023), we employ Flan-T5 (Chung et al., 2022) as our chosen LLM and created several LoRAs based on FLAN datasets. We conducted

extensive experiments across various tasks, including Translation, Natural Language Inference (NLI), Struct to Text, Closed-Book QA, and multiple subtasks within the Big-Bench Hard (BBH) (Ghazal et al., 2013) dataset. We train 800 iterations for each required composition of LoRAs with an initial learning rate of  $1e-5$ , batch size 12 and  $\alpha$  as 0.5.

**Compared Baselines.** We compared our MOLE with recently released state-of-the-art LoRA composition methods: LoRAhub (Han et al., 2023) and PEMs (Zhang et al., 2023).

**Main Results.** The corresponding experimental results are encapsulated in the Table 3. In summary, our MOLE surpasses state-of-the-art LoRA composition methods on five distinct datasets. Notably, on the BBH dataset, our MOLE achieves an average performance improvement of 3.8 over LoRAhub and outperforms PEMs by a notable margin of 9.0. Furthermore, in the realm of generation tasks, specifically in Translation and Struct to Text categories, MOLE consistently outshines its counterparts. In the Translation task set, it surpasses LoRAhub by an average margin of 1.5 and PEMs by 2.7. Correspondingly, within the Struct to Text task set, our model boasts an average performance superiority of 2.1 over LoRAhub and 2.6 over PEMs. These findings underscore the efficacy and versatility of our MOLE in handling language generation tasks.

## 5 ANALYSIS

### The effectiveness of gating balancing loss.

Figure 5 (a) and (b) illustrate how our  $\mathcal{L}_{\text{balance}}$  function mitigates the reduction in entropy rates within gating functions, leading to a more uniform composition weight distribution. The performance comparison between MOLE and MOLE  $_{w/o \mathcal{L}_{\text{balance}}}$  in Table 7 underscores the performance enhancement achieved with the inclusion of  $\mathcal{L}_{\text{balance}}$ . Additionally, we conducted an experiment wherein we solely increased the temperature  $\tau$  in Eq. 11, as an alternative to adding  $\mathcal{L}_{\text{balance}}$ . Results in Table 7 shows declining performance in MOLE variants MOLE $^{\tau_1}$ , MOLE $^{\tau_2}$ , MOLE $^{\tau_3}$  ( $\tau_1 < \tau_2 < \tau_3$ ) with increasing temperature. While temperature rise addresses gating imbalance, it restricts dynamic LoRA exploration in MOLE, leading to inferior outcomes.

**Further comparison with SOTA multi-concept generation methods.** In the absence of comparable LoRA composition methods in the V&L domain, we incorporated two leading multi-concept generation algorithms that do not utilize LoRA: Custom (Kumari et al., 2023) and Textual Inversion (Gal et al., 2022a), both of which emphasize full-parameter training for enhanced results. As presented in Table 2, MOLE outperforms Textual Inversion in both image and text alignment and excels over Custom in text alignment. Furthermore, it’s worth noting that our MoLE is more lightweight compared to these full-parameter training methods. These comparisons underscore the superior effectiveness of our MoLE relative to methods that involve extensive parameter tuning.

**Scale to a larger number of LoRAs.** We explore the performance as the number of LoRAs increases. In the NLP domain, experiments were conducted with varying numbers of LoRA (8, 24, 48, 128),

# Task	Metric	LoRAHub	PEMs	MoLE
<b>Translation</b>				
WMT '14 En→Fr	BLEU	<u>27.4</u>	25.6	<b>29.1</b>
WMT '14 Fr→En	BLEU	<u>29.4</u>	27.1	<b>31.3</b>
WMT '16 En→De	BLEU	24.6	<u>24.9</u>	<b>27.7</b>
WMT '16 De→En	BLEU	<b>29.9</b>	28.0	<u>29.1</u>
WMT '16 En→Ro	BLEU	<u>17.7</u>	15.2	<b>18.9</b>
WMT '16 Ro→En	BLEU	<u>23.5</u>	21.7	<b>25.1</b>
Average		<u>25.4</u>	24.2	<b>26.9</b>
<b>Struct to Text</b>				
CommonGen	Rouge-1	<u>53.7</u>	48.8	<b>55.1</b>
	Rouge-2	<b>23.1</b>	22.4	<u>23.1</u>
	Rouge-L	<u>49.7</u>	47.2	<b>53.9</b>
DART	Rouge-1	45.3	<u>46.2</u>	<b>48.8</b>
	Rouge-2	22.6	18.9	<b>23.5</b>
E2ENLG	Rouge-L	35.1	<b>37.6</b>	<u>36.0</u>
	Rouge-1	<u>41.1</u>	40.7	<b>42.0</b>
	Rouge-2	<u>26.3</u>	24.2	<b>29.0</b>
WebNLG	Rouge-L	38.8	<b>42.1</b>	<u>41.8</u>
	Rouge-1	<u>52.1</u>	52.0	<b>54.5</b>
	Rouge-2	23.9	<u>24.6</u>	<b>26.8</b>
Average	Rouge-L	45.2	<u>47.8</u>	<b>49.3</b>
<b>Closed-Book QA</b>				
ARC-c	EM	<u>51.7</u>	50.4	<b>52.9</b>
ARC-e	EM	<u>69.7</u>	65.7	<b>70.3</b>
NQ	EM	<u>17.3</u>	16.1	<b>23.5</b>
TQA	EM	<b>54.5</b>	53.9	<u>54.0</u>
Average		<u>48.3</u>	46.5	<b>50.2</b>
<b>Big-Bench Hard (BBH)</b>				
Boolean Expressions	EM	<u>55.1</u>	53.0	<b>57.3</b>
Causal Judgement	EM	<u>57.6</u>	51.1	<b>57.9</b>
Date Understanding	EM	<b>31.0</b>	29.3	<u>30.7</u>
Disambiguation	EM	46.6	<u>47.2</u>	<b>49.3</b>
Penguins in a Table	EM	<u>41.4</u>	39.8	<b>45.0</b>
Reasoning Objects	EM	<u>35.2</u>	<b>37.5</b>	33.7
Ruin Names	EM	<u>19.9</u>	19.3	<b>21.2</b>
Average		<u>38.4</u>	33.2	<b>42.2</b>
<b>Natural Language Inference (NLI)</b>				
ANLI-R1	EM	81.0	80.3	<b>82.7</b>
ANLI-R2	EM	<u>80.9</u>	80.2	<b>82.4</b>
ANLI-R3	EM	<u>77.4</u>	76.6	<b>78.9</b>
QNLI	EM	<u>77.6</u>	<u>78.0</u>	<b>78.1</b>
Average		<u>79.2</u>	78.8	<b>80.5</b>

Table 3: Evaluation results on Translation, Struct to Text, Closed-Book QA, NLI and BBH. The **best value** is in bold and the second-best value is underlined.



as detailed in Table 6. Our MOLE demonstrated optimal performance across these configurations, notably excelling with larger LoRA counts of 48 and 128, surpassing LoRAHub by **2.5** and **3.0**, respectively. Analysis revealed that LoRAHub’s optimization algorithm often zeroes out many LoRA weights in larger arrays, thus underutilizing the potential of all LoRA. Conversely, MOLE effectively overcomes this limitation. However, all methods, including MOLE, showed performance declines with an extremely large number of LoRA (128), highlighting a need for further research in this area. In the V&L domain, Table 10 shows experiments with increased composed LoRAs. While typical composition involve 3-4 visual concepts, our range was 3-6 to avoid ambiguity in outputs. Results indicate that MOLE consistently outperforms other LoRA composition models in text and image alignment as the number of LoRAs increases, underscoring its robustness and superior composition capabilities.

**Coarse-to-fine gating analysis.** To examine the impact of different granularity levels in gating functions, we delineated four levels in MOLE: matrix-wise (MOLE, gating at the parameter matrix level), layer-wise (MOLE), block-wise (MOLE), and network-wise (MOLE), abbreviated as m-MOLE, l-MOLE, b-MOLE, and n-MOLE respectively. Table 9 reveals that intermediate granularities, b-MOLE and l-MOLE, achieved the highest performance. In contrast, the coarsest level, n-MOLE, which involves minimal optimizable parameters (a single gating for the entire network), showed suboptimal outcomes. Additionally, the finest granularity, m-MOLE, underperformed, potentially due to its excessive control interfering with inherent relationships in LoRA parameters.

**Generalization to new datasets.** To further validate the effectiveness of our MOLE, we conducted generalization experiments. Specifically, all LoRA candidates and LoRA composition variants, including MOLE, PEMs and LoRAHub, were trained on NLI tasks (ANLI-R1, ANLI-R2, ANLI-R3, QNLI, and WNLI, among others). Subsequently, we evaluated these methods on the BBH dataset. As illustrated in Table 8, our MOLE achieves an average performance advantage of 2.4 over LoRAHub and 3.7 over PEMs, underscoring its superior generalization ability.

**Flexibility of MOLE.** As discussed in Section 2.1, a well-designed LoRA composition method should not only achieve effective LoRA composition but also retain the characteristics of individual LoRA. It should be versatile enough to function as a standalone LoRA generator, ensuring its practical applications are flexible and widespread. Figure 6 displays a comparison of the qualitative results for the retaining ability of several composition methods, we find that our MOLE can generate images that closely resemble the original features of the LoRA experts (e.g., dog ears, the color of the backpack), while other composition methods tend to produce confusion and loss of LoRA characteristics. Besides, as shown in Figure 1, we can also degrade MOLE by masking out the LoRA experts we do not wish to use, transforming it into a MOLE that merges fewer LoRAs without affecting the composition effect of the remaining LoRAs. As shown in Figure 8, our MOLE can achieve the same flexible LoRA composition as linear arithmetic composition method without altering the weights of MOLE, while reference tuning-based composition (Gu et al., 2023) can not accomplish.

**Hierarchical control analysis.** MOLE aims to achieve improved LoRA composition effects through finer-grained hierarchical control. As illustrated in the Figure 7, we visualize the weight distributions assigned by the gating functions learned by MOLE at different levels in both NLP and V&L domains. We observe that MOLE adaptively assigns weights to different LoRA experts at various layers. Consequently, finer-grained weight combination methods lead to superior results.

## 6 CONCLUSION AND LIMITATIONS

In this study, we introduce the Mixture of LoRA Experts (MOLE) as a versatile and dynamic approach for composing multiple trained LoRAs. The key innovation of MOLE lies in its learnable gating functions, which utilize the outputs of multiple LoRAs at each layer to determine composition weights. Our comprehensive evaluation in both the both NLP and V&L domains establishes that MOLE outperforms existing LoRA composition methods.

**Limitations.** As described in Section 5, when the number of LoRAs increases to a very large value (e.g., 128), despite our MOLE exhibiting superior performance, the performance of all LoRA composition methods, including our MOLE, tends to decrease. This suggests that our MOLE still faces challenges when performing large-scale LoRA composition. It also highlights the significance of researching better approaches for handling large-scale LoRA composition effectively.

## REFERENCES

- Shengnan An, Yifei Li, Zeqi Lin, Qian Liu, Bei Chen, Qiang Fu, Weizhu Chen, Nanning Zheng, and Jian-Guang Lou. Input-tuning: Adapting unfamiliar inputs to frozen pretrained models. *arXiv preprint arXiv:2203.03131*, 2022.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022a.
- Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022b.
- Ahmad Ghazal, Tilmann Rabl, Mingqing Hu, Francois Raab, Meikel Poess, Alain Crolotte, and Hans-Arno Jacobsen. Bigbench: Towards an industry standard benchmark for big data analytics. In *Proceedings of the 2013 ACM SIGMOD international conference on Management of data*, pp. 1197–1208, 2013.
- Yuchao Gu, Xintao Wang, Jay Zhangjie Wu, Yujun Shi, Yunpeng Chen, Zihan Fan, Wuyou Xiao, Rui Zhao, Shuning Chang, Weijia Wu, et al. Mix-of-show: Decentralized low-rank adaptation for multi-concept customization of diffusion models. *arXiv preprint arXiv:2305.18292*, 2023.
- Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris Metaxas, and Feng Yang. Svdiff: Compact parameter space for diffusion fine-tuning. *arXiv preprint arXiv:2303.11305*, 2023.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. Lorahub: Efficient cross-task generalization via dynamic lora composition. *arXiv preprint arXiv:2307.13269*, 2023.
- Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1931–1941, 2023.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial nli: A new benchmark for natural language understanding. *arXiv preprint arXiv:1910.14599*, 2019.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021a.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021b.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable questions for squad. *arXiv preprint arXiv:1806.03822*, 2018.

- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22500–22510, 2023.
- Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. VI-adapter: Parameter-efficient transfer learning for vision-and-language tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5227–5237, 2022.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman.  $p+$ : Extended textual conditioning in text-to-image generation. *arXiv preprint arXiv:2303.09522*, 2023.
- Yuan Xie, Shaohan Huang, Tianyu Chen, and Furu Wei. Moec: Mixture of expert clusters. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 13807–13815, 2023.
- Jinghan Zhang, Shiqi Chen, Junteng Liu, and Junxian He. Composing parameter-efficient modules with arithmetic operations. *arXiv preprint arXiv:2306.14870*, 2023.
- Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.

Table 4: The first motivation experiment in the NLP domain. NLA denotes normalized linear arithmetic composition (Eq. 2). The **best value** is in bold.

Model	ANLI-R1	ANLI-R2	ANLI-R3	QNLI	WNLI	Average
Single LoRA	<b>80.32</b>	<b>79.02</b>	75.92	<b>78.62</b>	<b>74.32</b>	<b>77.64</b>
NLA	79.32	78.88	<b>76.42</b>	78.06	69.98	76.53

Table 5: The second motivation experiment in the NLP domain. Full LoRA denotes the application of the complete set of LoRA parameters for inference, whereas x%-y% indicates the inference using LoRA parameters ranging from the top x% to the top y%. The **best value** is in bold.

	ANLI-R1	ANLI-R2	QNLI
Full LoRA	81.65	80.03	76.42
0%-20%	78.72	78.35	<b>78.14</b>
20%-40%	76.10	77.96	77.85
40%-60%	76.95	<b>81.47</b>	74.57
60%-80%	77.25	78.19	75.71
80%-100%	<b>82.59</b>	77.91	75.48

Table 6: NLP domain experimental results on the impact of exploring expand expert numbers on model performance. The result is the average EM on the Big-Bench Hard (BBH) dataset. NLA denotes normalized linear arithmetic composition (Eq. 2). The **best value** is in bold and the second-best value is indicated with an underline.

# Number of LoRA	NLA	LoRAHub	PEMs	MOLE
8	32.7	<u>33.9</u>	33.7	<b>36.6</b>
24	36.8	<u>37.1</u>	36.9	<b>38.7</b>
48	34.4	<u>36.9</u>	34.6	<b>39.4</b>
128	34.1	<u>35.5</u>	34.9	<b>38.5</b>
Average	34.5	<u>35.9</u>	35.0	<b>38.3</b>

Table 7: Experimental results on gating balance of MOLE. NLA denotes normalized linear arithmetic composition (Eq. 2). The **best value** is in bold.

# Model	ANLI-R1	ANLI-R2	ANLI-R3	QNLI	WNLI	Average
NLA	79.32	78.88	76.42	78.06	69.98	76.53
MOLE	<b>81.49</b>	<b>79.38</b>	<b>77.63</b>	<b>79.52</b>	<b>72.31</b>	<b>78.07</b>
MOLE w/o $\mathcal{L}_{\text{balance}}$	80.81	79.11	77.42	79.09	71.44	77.57
MOLE <sup>T1</sup>	80.52	79.27	77.30	79.11	71.07	77.45
MOLE <sup>T2</sup>	80.01	79.03	76.33	77.81	70.37	76.71
MOLE <sup>T3</sup>	78.50	79.20	76.07	78.02	70.00	76.35

Table 8: Evaluation results on generalization to new datasets. All lora candidates and LoRA merging variants are optimized on NLI tasks. The **best value** is in bold and the second-best value is indicated with an underline.

# Task	Metric	LoRAHub	PEMs	MOLE
<b>Big-Bench Hard (BBH)</b>				
Boolean Expressions	EM	45.3	<u>45.5</u>	<b>48.7</b>
Causal Judgement	EM	<u>51.3</u>	46.1	<b>52.4</b>
Date Understanding	EM	<b>27.5</b>	24.6	<u>26.6</u>
Disambiguation	EM	39.7	<u>42.4</u>	<b>43.8</b>
Penguins in a Table	EM	<u>35.3</u>	33.6	<b>39.0</b>
Reasoning about Colored Objects	EM	<u>32.2</u>	31.4	<b>34.7</b>
Average		<u>38.5</u>	37.2	<b>40.9</b>

Table 9: Coarse-to-fine gating comparison. The **best value** is in bold and the second-best value is indicated with an underline.

# Method	Text-alignment	Image-alignment		
		Concept 1	Concept 2	Concept 3
m-MoLE	0.731	0.719	0.714	0.747
l-MoLE	<u>0.760</u>	<u>0.727</u>	<u>0.731</u>	<b>0.757</b>
b-MoLE	<b>0.766</b>	0.726	<b>0.737</b>	<u>0.755</u>
n-MoLE	0.722	<b>0.739</b>	0.682	0.730

Table 10: Experimental results on the impact of exploring expand expert numbers on model performance. We evaluate each composition pair on 200 images generated using 5 prompts with 50 steps of DDPM sampler and scale=7.5. NLA denotes normalized linear arithmetic composition (Eq. 2). The best performance is in bold.

# Number of LoRA	Text-alignment			Average Image-alignment		
	NLA	SVDiff	MoLE	NLA	SVDiff	MoLE
3	0.678	0.728	<b>0.759</b>	0.694	0.719	<b>0.757</b>
4	0.681	0.717	<b>0.725</b>	0.712	0.721	<b>0.742</b>
5	0.652	0.723	<b>0.762</b>	0.682	0.708	<b>0.737</b>
6	0.698	0.709	<b>0.737</b>	0.703	0.701	<b>0.709</b>
Average	0.677	0.719	<b>0.746</b>	0.698	0.712	<b>0.736</b>

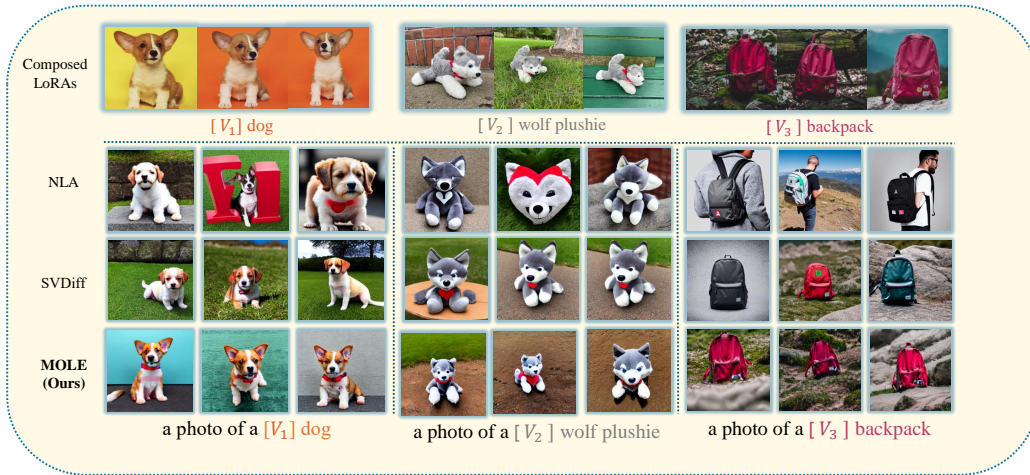


Figure 6: Qualitative result for retaining ability experiment. NLA denotes normalized linear arithmetic composition (Eq. 2). The first row displays the composed trained LoRAs. The second to the last row showcases the respective abilities of different composition methods to preserve the characteristics of each LoRA without altering the model.

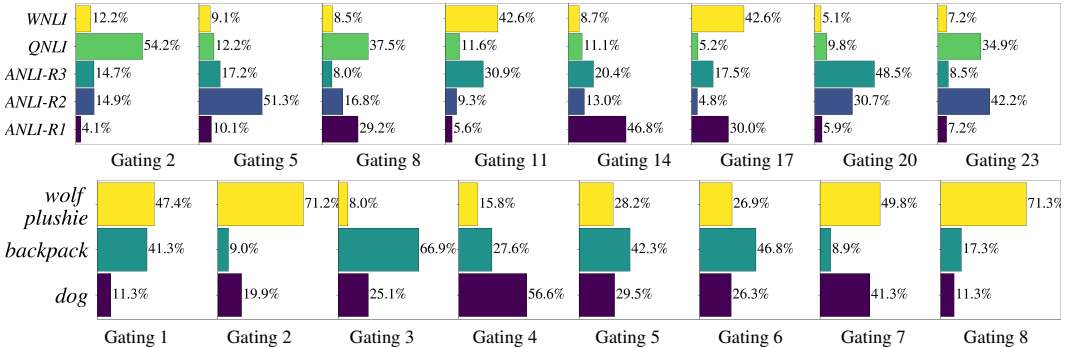


Figure 7: Visualization of the weights (%) predicted by each gating function (horizontal axis) for LoRA experts (vertical axis) during inference. The top row corresponds to experiments in the NLP domain, while the bottom row pertains to experiments in the V&L domain.



Figure 8: Visualization for different inference modes of MOLE. MOLE has two inference modes: In the first mode (the first line), MOLE can use all the LoRA experts and allocate weights for each LoRA, preserving their individual characteristics. In the second mode (the second and third lines), we can manually mask some unwanted LoRAs without changing the gating weights. It can recalculate and distribute weights proportionally. These two modes enable MOLE to adapt to different scenarios, providing a versatile and flexible approach for effective LoRA composition.



Figure 9: Visualization of multiple LoRA composition results on V&L domain. NLA denotes normalized linear arithmetic composition (Eq. 2). Our MOLE has higher visual similarity with the personal cat and dog images while following the text condition better, e.g., SVDiff is unable to fully recover all the characteristics of LoRA (in the second line, the appearance of the dog is completely altered, and in the first line, two cats are present but the dog is missing). Moreover, SVDiff and NLA struggles to generate images that match the text condition effectively (e.g., it might add sunglasses to both dogs and cats in response to conditions mentioning “dog” and “cat”).



Figure 10: Visualization of multiple LoRA composition results on V&L domain. NLA denotes normalized linear arithmetic composition (Eq. 2). Our model consistently produces results that better align with the prompt descriptions. The outputs from our model consistently contain all three visual concepts that need to be combined. In contrast, SVDiff and NLA often exhibit issues such as concept confusion (e.g., in the third row of NLA, where features of both the cat and dog are confused) and concept omission (e.g., in the second row of SVDiff, where the concept of the dog is missing, and in the first row, where the concept of the cat is missing).





Figure 11: Visualization of multiple LoRA composition results on V&L domain. NLA denotes normalized linear arithmetic composition (Eq. 2). Our model consistently produces results that better align with the prompt descriptions. The outputs from our model consistently contain all three visual concept features that need to be combined. In contrast, SVDiff and NLA often exhibit issues such as concept omission (e.g., in the first row of NLA, where the concepts of the cat and sunglasses are missing, and in the first row of SVDiff, where the concept of sunglasses is missing). Additionally, our output results better match the original visual concept features. For example, the shell of the turtle is green, whereas SVDiff and NLA generate shells in pink and brown colors.