# LLaVA Steering: Visual Instruction Tuning with 500x Fewer Parameters through Modality Linear Representation-Steering

Jinhe Bi<sup>1,2\*</sup> Yujun Wang<sup>1\*</sup> Haokun Chen<sup>1</sup> Xun Xiao<sup>2†</sup> Artur Hecker<sup>2</sup>

Volker Tresp<sup>1,3</sup> Yunpu Ma<sup>1,3†</sup>

<sup>1</sup> Ludwig Maximilian University of Munich <sup>2</sup> Munich Research Center, Huawei Technologies <sup>3</sup> Munich Center for Machine Learning

## Abstract

Multimodal Large Language Models (MLLMs) enhance visual tasks by integrating visual representations into large language models (LLMs). The textual modality, inherited from LLMs, enables instruction following and in-context learning, while the visual modality boosts downstream task performance through rich semantic content, spatial information, and grounding capabilities. These modalities work synergistically across various visual tasks. Our research reveals a persistent imbalance between these modalities, with text often dominating output generation during visual instruction tuning, regardless of using full or parameter-efficient fine-tuning (PEFT). We found that re-balancing these modalities can significantly reduce trainable parameters, inspiring further optimization of visual instruction tuning. To this end, we introduce Modality Linear Representation-Steering (MoReS), which re-balances intrinsic modalities by steering visual representations through linear transformations in the visual subspace across each model layer. We validated our approach by developing LLaVA Steering, a suite of models using MoReS. Results show that LLaVA Steering requires, on average, 500 times fewer trainable parameters than LoRA while maintaining comparable performance across three visual benchmarks and eight visual question-answering tasks. Finally, we introduce the LLaVA Steering Factory, a platform that enables rapid customization of MLLMs with a component-based architecture, seamlessly integrating state-of-theart models and evaluating intrinsic modality imbalance. This open-source project facilitates a deeper understanding of MLLMs within the research community. Code is available at https://github.com/bibisbar/LLaVA-Steering.

## 1 Introduction

Recent advancements in Multimodal Large Language Models (MLLMs) (Liu et al., 2024b; Xue et al., 2024; Zhou et al., 2024a; Bi et al., 2025a) have demonstrated impressive capabilities across a variety of visual downstream tasks (zhou changshi et al., 2025; Wang et al., 2025a; Huang et al., 2025a,c; Cui et al., 2025; Li et al., 2024c,a; Yue et al., 2025; Zhao et al., 2025; Zhang et al., 2023a). These models integrate visual representations from pretrained vision encoders via various connectors (Liu et al., 2024a; Li et al., 2023a; Alayrac et al., 2022) into LLMs, leveraging the latter's sophisticated reasoning abilities (Zhang et al., 2024a; Abdin et al., 2024; Zheng et al., 2023a).

To better integrate visual representations into LLMs, the most popular MLLMs adopt a two-stage training paradigm: pretraining followed by visual instruction tuning. In the pretraining stage, a connector is employed to project visual representations into the textual representation space. We define these two modalities—text and vision—as intrinsic to MLLMs, each carrying rich semantic information that serves as the foundation for further visual instruction tuning on downstream tasks such as image understanding (Sidorov et al., 2017a; Lu et al., 2022; Hudson and Manning, 2019), and instruction following (Liu et al., 2023a).

In the visual instruction tuning stage, due to its high computational cost, researchers have pursued two primary strategies. One approach focuses on refining data selection methodologies (Liu et al., 2024e; McKinzie et al., 2024) to reduce redundancy and optimize the training dataset, though this process remains expensive and time-consuming. A more common strategy goes to employ Parameter-Efficient Fine-Tuning (PEFT) methods (Huang et al., 2025b; Zhang et al., 2025), such as LoRA (Hu et al., 2021), aiming to reduce the number

<sup>\*</sup>These authors contributed equally to this work. Email contact: *bijinhe@outlook.com*.

<sup>&</sup>lt;sup>†</sup>Corresponding authors: *cognitive.yunpu@gmail.com*, *drxiaoxun@gmail.com* 



Figure 1: Left: Attention score distributions across layers for three MLLM fine-tuning methods (Full, LoRA, and MoReS), sampled from 100 instances each. Green represents visual representations, while grey indicates other (primarily textual) representations. Full fine-tuning and LoRA show strong reliance on textual representations across most layers. In contrast, the proposed MoReS method demonstrates significantly improved visual representation utilization, particularly in the middle and lower layers, addressing the intrinsic modality imbalance in MLLMs. **Right:** Average visual attention score distribution versus model size for different MLLM fine-tuning methods. The plot suggests that methods achieving better balanced intrinsic modality tend to require fewer trainable parameters.

of trainable parameters, thereby making visual instruction tuning more computationally feasible (Liu et al., 2024a; Zhou et al., 2024a). However, even with PEFT methods like LoRA, large-scale MLLMs remain prohibitively expensive to fine-tuning.

This raises a critical question: is there any further possibility to reduce more trainable parameters so that the visual instruction tuning can be further improved? Our research offers a novel viewpoint by focusing on the intrinsic modality imbalance within MLLMs. A closer analysis uncovers an imbalance in output attention computation (Chen et al., 2024a), where textual information tends to dominate the attention distribution during output generation. Specifically, we investigate this issue by analyzing attention score distributions, which evaluates the balance between text and visual modalities. As shown in Figure 1, visual representations are significantly underutilized during visual instruction tuning. More importantly, our analysis reveals that achieving a better balance between these modalities can substantially reduce the number of trainable parameters required for fine-tuning. Hereby we suppose that intrinsic modality rebalance is the Midas touch to unlock further reductions in the number of trainable parameters.

To address this challenge, we introduce Modality Linear Representation-Steering (MoReS) to optimize visual instruction tuning, significantly reducing the number of trainable parameters while maintaining equivalent performance. Unlike full finetuning, which modifies the entire model, or other popular PEFT methods such as LoRA (Hu et al., 2021), OFT (Qiu et al., 2023), Adapter (Houlsby et al., 2019), and IA3 (Liu et al., 2022a), MoReS focuses solely on steering the visual representations. Specifically, our approach freezes the entire LLM during visual instruction tuning to preserve its capabilities in the textual modality. Instead of fine-tuning the full model, we introduce a simple linear transformation to steer visual representations in each layer. This transformation operates within a subspace after downsampling, where visual representations encode rich semantic information in a compressed linear subspace (Zhu et al., 2024; Shimomoto et al., 2022; Yao et al., 2015). By continuously steering visual representations across layers, MoReS effectively controls the output generation process, yielding greater attention inclined to visual modality.

To validate the efficacy of our proposed MoReS method, we integrated it into MLLMs of varying scales (3B, 7B, and 13B parameters) during visual instruction tuning, following the LLaVA 1.5 (Liu et al., 2024a) training recipe. The resulting models, collectively termed LLaVA Steering, achieved competitive performance across three visual benchmarks and six visual question-answering tasks, while requiring 287 to 1,150 times fewer trainable parameters than LoRA, depending on the specific training setup.

In our experiments, we observed the need for a comprehensive framework to systematically analyze and compare various model architectures and training strategies in MLLMs. The wide range of

design choices and techniques makes it difficult to standardize and understand the interplay between these components. Evaluating each method across different open-source models is time-consuming and lacks consistency due to implementation differences, requiring extensive data preprocessing and careful alignment between architectures and training recipes. To address this issue, we developed the LLaVA Steering Factory, a flexible framework that reimplements mainstream vision encoders, multiscale LLMs, and diverse connectors, while offering customizable training configurations across a variety of downstream tasks. This framework simplifies pretraining and visual instruction tuning, minimizing the coding effort. Additionally, we have integrated our attention score distribution analysis into the LLaVA Steering Factory, providing a valuable tool to the research community for further studying intrinsic modality imbalance in MLLMs. Our work makes the following key contributions to the field of MLLMs:

- 1. First of all, we propose Modality Linear Representation-Steering (MoReS), a novel method that addresses intrinsic modality imbalance in MLLMs by steering visual representations through linear transformations within the visual subspace, effectively mitigating the issue of text modality dominating visual modality.
- 2. In addition, we present LLaVA Steering, where with different sizes (3B/7B/13B), three real-world LLaVA MLLMs consisting of different model components are composed by integrating the proposed MoReS method into visual instruction tuning. LLaVA Steering models based on MoReS method achieve comparable performance across three visual benchmarks and six visual question-answering tasks, while requiring 287 to 1,150 times fewer trainable parameters.
- 3. Last but not least, we develop the LLaVA Steering Factory, a flexible framework designed to streamline the development and evaluation of MLLMs with minimal coding effort. It offers customizable training configurations across diverse tasks and incorporates tools such as attention score analysis, facilitating systematic comparisons and providing deeper insights into intrinsic modality imbalance.

# 2 Related Work

**Integrating Visual Representation into LLMs:** Existing approaches for integrating visual representations into LLMs broadly fall into three categories: (1) Cross-attention architectures (e.g., Flamingo (Awadalla et al., 2023), IDEFICS (Laurençon et al., 2023)) that inject image features through adapter layers while keeping LLM weights frozen; (2) Decoder-only architectures like LLaVA (Liu et al., 2024b) and Qwen-VL (Bai et al., 2023) that train visual projectors during pretraining and often unfreeze LLMs during fine-tuning; and (3) Visionencoder-free methods (Chen et al., 2024b; Diao et al., 2024) that process raw pixels directly. Hybrid approaches like NVLM (Dai et al., 2024) combine elements of these paradigms. While effective, these methods incur substantial computational costs during visual instruction tuning due to large-scale multimodal alignment requirements.

**Visual Instruction Tuning:** Fine tuning of multimodal large language models (MLLMs) for downstream tasks has gained considerable attention, but remains computationally expensive due to largescale visual instruction datasets and model sizes (Wang et al., 2022). To tackle this challenge, recent advancements have introduced parameter-efficient fine-tuning (PEFT) methods (Houlsby et al., 2019; Li and Liang, 2021), such as LoRA (Hu et al., 2021), enabling more efficient visual instruction tuning.

However, many of these PEFT methods primarily focus on optimizing weights but ignore the intrinsic representation imbalance during visual instruction tuning, thus cannot further reduce the required trainable parameters. This means to look for other novel approaches that can improve the efficiency and effectiveness of visual instruction tuning.

**Representation Steering:** Recent studies (Singh et al., 2024; Avitan et al., 2024; Li et al., 2024b; Subramani et al., 2022) have demonstrated that the representations induced by pre-trained language models (LMs) encode rich semantic structures. Steering operations within this representation space have shown to be effective in controlling model behavior. Unlike neuron-based or circuit-based approaches, representation steering manipulates the representations themselves, providing a clearer mechanism for understanding and controlling the behavior of MLLMs and LLMs. For example, (Zou et al., 2023) explores representation engineering to modify neural network behavior, shifting the



Figure 2: Layer-wise Modality Attention Ratio (LMAR) comparison across training methods, including Full finetuning, LoRA, Adapter, IA3, and our MoReS. Our MoReS method (red line) consistently demonstrates the highest LMAR across most layers, with a notable spike in the final layers. Compared with full fine-tuning and mainstream PEFT methods, our MoReS needs the least parameters during visual instruction tuning while achieving superior modality balance.

focus from neuron-level adjustments to transformations within the representation space. Similarly, (Wu et al., 2024a) applies scaling and biasing operations to alter intermediate representations. Furthermore, (Wu et al., 2024b) introduces a family of representation-tuning methods that allows for interpretable interventions within linear subspaces. In this work, we leverage the concept of representation steering to introduce a novel approach, MoReS, which enhances attention to visual representations, thereby demonstrating superior parameter efficiency compared to baseline PEFT methods (Hu et al., 2021; Houlsby et al., 2019; Liu et al., 2022a; Qiu et al., 2023).

# **3** Intrinsic Modality Imbalance

This section explores how the two intrinsic modalities—text and vision—are imbalanced during output generation across each layer in MLLMs, as reflected in the attention score distribution. Furthermore, we demonstrate that addressing this modality imbalance effectively during visual instruction tuning can guide the design of methods that require fewer trainable parameters.

We begin with calculating the attention score distribution across both modalities in each layer, as derived from the generated output. In auto-regressive decoding, which underpins decoder-only MLLMs, output tokens are generated sequentially, conditioned on preceding tokens. The probability distribution over the output sequence  $\hat{y}$  is formalized as:

$$p(\hat{y}) = \prod_{i=1}^{L} p(\hat{y}_i | \hat{y}_{< i}, R_{\text{text}}, R_{\text{image}}, R_{\text{sys}}) \quad (1)$$

where  $\hat{y}_i$  represents the *i*-th output token,  $\hat{y}_{<i}$  denotes the preceding tokens,  $R_{\text{text}}$  is the textual representation,  $R_{\text{image}}$  is the visual input representation,  $R_{\text{sys}}$  accounts for system-level contextual information, and L is the output sequence length.

To quantify modality representation imbalance, we calculate the sum of attention scores allocated to visual representations across all layers in MLLMs. Figure 1 illustrates this imbalance across full fine-tuning, LoRA, and our proposed MoReS method. The results indicate that textual representations often dominate the output generation process in both full fine-tuning and LoRA.

Further examination of this imbalance across multiple PEFT methods reveals an intriguing trend: methods that make better use of visual representations tend to require fewer trainable parameters during visual instruction tuning.

To validate this observation, we introduce the Layer-wise Modality Attention Ratio (LMAR), formulated as:

$$LMAR_{l} = \frac{1}{N} \sum_{i=1}^{N} \frac{\alpha_{l}^{\text{image},i}}{\alpha_{l}^{\text{text},i}} , \qquad (2)$$

where l denotes the layer index, N is the total number of samples, and  $\alpha_l^{\text{image},i}$  and  $\alpha_l^{\text{text},i}$  are the mean attention scores allocated to visual and textual tokens, respectively, in layer l for the *i*-th sample. LMAR thus provides a robust measure of the attention distribution between modalities, averaged over multiple samples to capture general trends in modality representation across layers. This value reflects balanced per-token attention allocation between visual and textual modalities, rather than total attention mass. Since visual inputs typically consist of hundreds of tokens (e.g., 576 for a  $24 \times 24$  patch grid), while textual inputs often include only dozens of tokens, simply summing attention weights across modalities can obscure imbalance. By normalizing attention at the per-token level, LMAR more accurately captures whether each modality is being fairly attended to, regardless of token count disparity. Therefore, an LMAR close to 1.0 implies that, on average, each



Figure 3: Schematic Overview of Modality Linear Representation-Steering (MoReS): Left: The architectural diagram depicts the integration of textual and visual tokens through transformer layers, leading to output token generation. **Right:** The mathematical formulation of MoReS illustrates the steering of visual representations within a subspace, highlighting its impact on output generation. During visual instruction tuning, the parameters of the LLM remain frozen, allowing only the parameters associated with the linear transformation in the steering mechanism to be trainable. With MoReS, the distribution of attention scores becomes more balanced, achieving intrinsic modality balance.

visual token receives comparable attention to each text token, which we interpret as a strong signal of modality balance within the model.

In our experiments comparing various existing PEFT methods and full fine-tuning, IA3 (Liu et al., 2022a) consistently achieves the highest average LMAR score across all layers while requiring the fewest trainable parameters. IA3's superior performance can be attributed to its unique design, which introduces task-specific rescaling vectors that directly modulate key components of the Transformer architecture, such as the keys, values, and feed-forward layers. Unlike methods that introduce complex adapters or fine-tune all parameters, IA3 optimizes a small but crucial set of parameters responsible for attention and representation learning. By applying element-wise scaling to the attention mechanisms, IA3 effectively re-balances the attention distribution across two intrinsic modalities. This design is particularly beneficial during visual instruction tuning, as it allows the model to dynamically reallocate more attention to visual representations without requiring many trainable parameters. The identified relationship inspires that if the intrinsic modality imbalance can be addressed, the required number of trainable parameters can be

potentially reduced further during visual instruction tuning. This offers a new direction for future improvements in PEFT methods for MLLMs.

## 4 MoReS Method

Based on insights gained from intrinsic modality imbalance, we introduce Modality Linear Representation-Steering (**MoReS**) as a novel method for visual instruction tuning which can rebalance visual and textual representations and achieve comparable performance with fewer trainable parameters.

Our approach is grounded in the linear subspace hypothesis, originally proposed by Bolukbasi et al. (2016), which suggests that information pertaining to a specific concept is encoded within a linear subspace in a model's representation space. This hypothesis has been rigorously validated across numerous domains, including language understanding and interpretability (Lasri et al., 2022; Nanda et al., 2023; Amini et al., 2023; Wu et al., 2024c). Building upon the intervention mechanisms described in Geiger et al. (2024) and Guerner et al. (2023), we introduce a simple linear transformation that steers visual representations within subspace while keeping the entire LLM frozen during visual instruction tuning. This approach ensures that the language model's existing capabilities are preserved, while continuously guiding the MLLM to better leverage the underutilized visual modality. By steering visual representations across each layer, MoReS effectively rebalances the intrinsic modality and influences the output generation process. Figure 3 provides an illustration of the overall concept and architecture behind MoReS.

Formally, MoReS method can be formulated as follows: Let  $\mathcal{H} = \{h_i\}_{i=1}^N \subset \mathbb{R}^D$  denote the set of visual representations in the original high-dimensional space. We define our steering function MoReS as:

$$MoReS(h) = W_{up} \cdot \phi(h) \tag{3}$$

where  $h \in \mathbb{R}^D$  is an input visual representation,  $\phi : \mathbb{R}^D \to \mathbb{R}^d$  is a linear transformation function that steers h into a lower-dimensional subspace  $\mathbb{R}^d$  (d < D), and  $W_{up} \in \mathbb{R}^{D \times d}$  is an upsampling matrix that projects from  $\mathbb{R}^d$  back to  $\mathbb{R}^D$ . The steering function  $\phi$  is defined as:

$$\phi(h) = \text{Linear}(h) - W_{\text{down}}h \tag{4}$$

where  $W_{\text{down}} \in \mathbb{R}^{d \times D}$  is a downsampling matrix. To preserve representational fidelity and ensure approximate invertibility, we impose the constraint  $W_{\text{down}}W_{\text{up}}^T = I_D$ .

This residual form disentangles the learnable modulation from the identity projection. Specifically,  $W_{\rm down}h$  projects the original representation into a lower-dimensional subspace, Linear(h) learns task-specific modulation, and the subtraction isolates the residual semantic shift. This design ensures a controlled, low-rank steering of the original representation, aligning with MoReS's goal of minimally altering the pretrained structure. Notably, this steering method can dynamically be applied to specific visual tokens. Further exploration of the impact of different steered token ratios is discussed in Section 5.7. In Section A.4, we further provide theoretical justification that elucidates how MoReS effectively rebalances the intrinsic modalities while continuously controlling output generation. Additionally, we provide a preliminary estimation of the trainable parameters involved during visual instruction tuning. In the following sections, we first compose real-world MLLMs (i.e., LLaVA Steering) with three different scales and integrate the proposed MoReS method. Based on the composed real-world models, we then evaluate how our MoReS method performs within the



Figure 4: Comparison of parameter count vs. performance for MoReS and other PEFT methods across four benchmarks.

composed models across several popular and prestigious datasets.

# **5** Experiments

We incorporate MoReS into each layer of the LLM during visual instruction tuning, developing LLaVA Steering (3B/7B/13B) based on the training recipe outlined in (Liu et al., 2024a). During visual instruction tuning on the LLaVA-665k dataset, we apply MoReS to a specific ratio of the total visual tokens, specifically using it on only 1% of the tokens. Further details about the model architectures and baseline training methods are provided in Appendix A.1.

#### 5.1 Multi-Task Supervised Fine-tuning

To assess the generality of our method, we compare it with the baselines using the LLaVA-665K multitask mixed visual instruction dataset (Liu et al., 2024a). Our evaluation covers several benchmarks, including VQAv2, GQA, VizWiz, ScienceQA, TextVQA, MM-Vet, POPE, and MMMU, to evaluate the performance across a range of tasks, from visual perception to multimodal reasoning. Further details can be found in Appendix A.2.

Following (Zhou et al., 2024b), we define ScienceQA as an unseen task, while VQAv2, GQA, and VizWiz are categorized as seen tasks in LLaVA-665k. To provide a comprehensive evaluation of our MoReS capabilities, we design three configurations: MoReS-Base, MoReS-Large, and MoReS-Huge, each based on different ranks. We present the results in Table 1, where our MoReS method achieves the highest scores on POPE (88.2) and MMMU (35.8), as well as the second-best performance on ScienceQA (71.9) and MM-Vet (33.3). Notably, MoReS accomplishes these results with 287 to 1150 times fewer trainable parameters compared to LoRA. The scatter plots in Figure 4 further illustrate that MoReS variants (highlighted in red) consistently achieve Pareto-optimal performance, offering an ideal balance between model size and effectiveness.

# 5.2 Task-Specific Fine-tuning

We evaluate the task-specific fine-tuning capabilities of our MoReS method in comparison to other tuning methods on multiple visual question answering datasets: (1) ScienceQA-Image (Lu et al., 2022), (2) VizWiz (Gurari et al., 2018), and (3) IconQA-txt and IconQA-blank (Lu et al., 2021). We present the results in Table 2, showing that MoReS achieves 1200 times fewer trainable parameters compared to LoRA and 3 times fewer than the previous best, IA3, while maintaining comparable performance or an acceptable decline of less than 3%. These results show that MoReS can succeed at Task-Specific Fine-tuning, even unseen tasks during its multitask visual instruction tuning stage.

# 5.3 Multi-scale Data Fine-tuning

During visual instruction tuning, the scale of specific task datasets can vary significantly. To gain a comprehensive understanding of our method compared to other training approaches, we follow the methodology of (Chen et al., 2022) and randomly sample 1K, 5K, and 10K data points from each dataset, defining these as small-scale, mediumscale, and large-scale tasks, respectively. Given the limited resources available, we choose MoReS-L for fine-tuning.

Table 3 demonstrates that MoReS exhibits strong capabilities across all scales. Notably, in small-scale tasks, MoReS outperforms full fine-tuning performance while using only 575 times fewer parameters than LoRA and 8,475 fewer than full fine-tuning. In contrast, methods like OFT and IA3 fail to surpass full fine-tuning despite utilizing significantly more parameters. This result underscores the practicality of MoReS in real-world scenarios where data collection can be challenging, suggesting that MoReS is suitable for multi-scale visual instruction tuning.

# 5.4 Text-only Tasks

MoReS preserves 100% of the pre-trained world knowledge in the LLM by neither modifying its parameters nor interfering with textual token inference. This design allows MoReS to excel in understanding both visual and textual information. Unlike many existing methods, which often alter model weights and risk degrading pre-trained knowledge (Zhang et al., 2024c), MoReS employs a representation-steering approach to selectively enhance the performance of the visual modality. Table 4 clearly demonstrate that MoReS excels in text-only tasks, further emphasizing its ability to retain and effectively leverage the inherent world knowledge stored in LLMs. This capability showcases MoReS' generalizability not only for multimodal tasks but also for text-dominant tasks.

# 5.5 Hallucination Mitigation

Hallucination remains a critical challenge in MLLMs, largely due to their strong linguistic bias, which can overshadow visual information and lead to outputs misaligned with the provided visual context. MoReS significantly outperforms existing tuning approaches in mitigating hallucinations, as demonstrated through evaluations on two widely recognized benchmarks: POPE and Hallucination-Bench. Key metrics include *Acc*, *Hard Acc*, *Figure Acc*, and *Question Acc*. Further details can be found in Appendix A.3.

Table 6 highlights the robustness of MoReS in reducing hallucination and enhancing the balance between linguistic and visual information in MLLMs.

# 5.6 Dynamic Steering Ratio

To explore the adaptability of MoReS across tasks of varying complexity, we introduce a dynamic steering ratio mechanism that allows the proportion of modulated visual tokens to vary during training. For each downstream task, we predefine a range of steering ratios and enable the model to dynamically select among them.

As shown in Table 9, performance trends differ across tasks. On SQA and IconQA, increasing the steering ratio beyond 25% leads to performance degradation, likely due to oversteering that disturbs the pretrained visual subspace. In contrast, VizWiz—a visually complex and low-quality dataset—benefits from higher steering ratios, suggesting that more intensive modulation is needed to adapt to noisy or ambiguous visual inputs. These

Method	TP*	VQAv2	GQA	TextVQA	SciQA-IMG	POPE	MM-Vet	MMMU	Avg
FT	2.78B	79.2	61.6	57.4	71.9	87.2	35.0	38.2	61.5
Adapter	83M	77.1	58.9	53.5	68.1	86.7	29.4	34.2	58.2
LoRA	188.7M	77.6	59.7	53.8	71.6	87.9	33.3	35.6	59.9
OFT	39.3M	75.1	55.3	52.9	69.1	87.6	31.0	35.6	58.3
IA3	0.49M	74.5	52.1	49.3	72.2	86.9	30.9	34.3	57.1
MoReS-B	0.164M	74.1	52.1	48.5	70.0	87.6	30.3	35.3	56.9
MoReS-L	0.328M	74.0	51.6	49.3	71.6	87.2	33.3	34.4	57.3
MoReS-H	0.655M	74.2	51.8	48.3	71.9	88.2	31.1	35.8	57.4

Table 1: Experimental results of Multi-Task Supervised Fine-tuning. TP\* denotes the number of trainable parameters within the LLM. All models share the same training recipe for the vision encoder and connector for fair comparison.

Model	Method	TP*	SciQA-IMG	VizWiz	IconQA-text	IconQA-blank
	Adapter	83M	92.3	62.9	93.5	95.8
T.T 37A	LoRA	188.7M	93.9	61.6	93.9	96.5
LLavA Steering 2P	OFT	39.3M	86.3	42.0	87.8	42.0
Steering-5D	IA3	0.492M	90.2	58.4	84.5	94.7
	MoReS-B	0.164M	89.7	59.2	84.0	94.2
	Adapter	201.3M	82.7	59.7	72.1	71.6
T.T 37A	LoRA	319.8M	87.6	60.6	77.7	70.2
LLavA Steering 7P	OFT	100.7M	78.3	55.1	19.4	22.7
Steering-7B	IA3	0.614M	83.8	54.3	65.1	70.4
	MoReS-B	0.262M	83.6	54.2	64.2	70.2
	Adapter	314.6M	87.9	61.4	78.2	73.0
T.T 37A	LoRA	500.7M	92.1	62.0	80.2	73.2
LLavA Steering 12P	OFT	196.6M	82.7	59.5	3.4	22.3
Steering-15B	IA3	0.963M	90.5	54.6	73.8	71.7
	MoReS-B	0.410M	89.5	54.3	74.9	71.5

Table 2: Task-specific fine-tuning results for different LLaVA Steering scales. TP\* denotes trainable parameters within the LLM.

results indicate that dynamically adjusting the steering ratio based on task characteristics is a promising direction. It allows MoReS to maintain minimal intervention when possible, while scaling up adaptation for more challenging scenarios.

#### 5.7 Ablation Studies

To gain deeper insights into our MoReS method, we conduct ablation studies focusing on its subspace choice and steered visual token ratio. We use LLaVA Steering-3B model as our baseline for comparison. Table 7 and 8 summarize the results of two types of ablations.

First, concerning the choice of subspace rank, we found that a rank of 1 achieves the highest average performance of 81.8 across four visual tasks while also requiring the fewest parameters, specifically 0.164M. Second, regarding the steered visual token ratio, we varied this parameter from 100% (dense steering) to 1% (sparse steering). The results indicate that a ratio of 1% is optimal, yielding the best or near-optimal performance on four benchmarks while also significantly reducing inference

Scale	Method	TP*	SciQA-IMG	VizWiz	IconQA
Small	FT	2.78B	33.8	51.2	68.1
	Adapter	83M	81.0	57.4	72.4
	LoRA	188.7M	84.0	58.5	74.2
	OFT	39.3M	79.2	43.2	35.9
	IA3	0.492M	79.9	50.5	73.0
	MoReS-L	0.328M	78.2	55.0	69.7
Medium	FT	2.78B	78.2	58.9	92.2
	Adapter	83M	92.1	60.6	93.2
	LoRA	188.7M	92.9	60.5	92.7
	OFT	39.3M	86.4	44.4	45.5
	IA3	0.492M	91.9	57.1	90.6
	MoReS-L	0.328M	92.1	56.6	89.9
Large	FT	2.78B	88.9	59.4	95.7
	Adapter	83M	92.4	61.3	95.2
	LoRA	188.7M	93.9	61.8	96.0
	OFT	39.3M	86.4	44.2	43.7
	IA3	0.492M	90.3	57.9	93.8
	MoReS-L	0.328M	89.8	57.7	93.5

Table 3: Results of multi-scale tasks. TP\* denotes trainable parameters within the LLM.

Task	LoRA	Adapter	OFT	IA3	MoReS (Ours)
HellaSwag	70.5	66.4	69.1	71.8	71.9
MMLU	55.3	52.9	54.7	56.8	57.0

Table 4: Performance comparison of PEFT methods on text-only tasks.

overhead due to its sparse steering approach.

#### 5.8 Discussion on Optimization

While MoReS is proposed as a standalone representation steering strategy, its design is inherently compatible with other model compression techniques such as sparse training and quantization. These methods have shown great promise in reducing the memory and computational footprint of large-scale models. However, MoReS is fundamentally distinct in its design philosophy. It preserves the full parameter set of the pretrained language model by freezing all original weights

Factory	Multi-scale LLMs I	Diverse Vision Encoders	PEFTs	Text-only Tasks	s Multimodal Tasks	Computational Optimization	Multiple Training Strategies
TinyLLaVA	×	1	X	×	1	×	✓
Prismatic	1	1	x	×	1	×	×
LLaVA Steering (Ours)	11	1	1	1	11	1	✓

Table 5: Comparison of functionality across different factories.

	Matria	Esti	LoDA	Adaptan	OFT	142	MaDa
	Metric	ruii	LOKA	Adapter	OFI	IAS	workes
POPE	Acc↑	87.2	86.7	87.9	85.1	86.9	88.2
HallucinationBench	Hard Acc↑	37.4	34.6	36.2	33.9	39.3	42.6
HallucinationBench	Figure Acc↑	18.5	16.7	18.2	14.1	18.5	19.4
HallucinationBench	Question Acc $\uparrow$	44.4	43.0	44.8	36.2	45.0	46.1

Table 6: Comparison of MoReS against other tuning methods on POPE and HallucinationBench benchmarks.

Subspace	Rank TP*	SciQA-IMG	VizWiz	IconQA-txt	IconQA-blank	. Avg
1	0.164M	89.6	59.2	84.0	94.2	81.8
2	0.328M	89.7	59.2	83.9	94.0	81.7
4	0.655M	89.5	58.7	83.8	94.1	81.5
8	1.340M	89.6	58.9	83.7	93.9	81.5

Table 7: Results of the subspace rank choice. The grey shading indicates the best results and our selected parameters.

and steering only a small number of visual token representations through lightweight subspace projections. This stands in contrast to sparsity- or quantization-based methods, which typically involve structural pruning or reduced-precision representations that may alter the internal behavior of the model. To maintain a clear analysis of modality rebalancing via representation steering, we intentionally do not integrate these additional techniques in this work. Nevertheless, due to the modularity and non-intrusive nature of MoReS (e.g., the use of linear transformation layers), such combinations are feasible. Early-stage experiments on integrating MoReS with sparsity and quantization are underway, and we consider this a promising direction for future research in building efficient and compact multimodal systems.

# 6 LLaVA Steering Factory

The LLaVA Steering Factory addresses the need for a comprehensive framework to systematically analyze and compare various MLLM architectures and training strategies. Standardizing the evaluation of these models is challenging due to implementation differences and diverse design choices. The LLaVA Steering Factory offers standardized training pipelines, flexible data preprocessing, and customizable model configurations. It supports mainstream LLMs, vision encoders, and PEFT methods, including our MoReS technique, and integrates in-

teered Visual Token	Ratio SciQA-IMG	VizWiz	IconQA-txt	IconQA-blank
1%	89.7	59.2	84.0	94.1
25%	89.9	59.0	80.2	93.8
50%	88.9	59.0	79.8	92.6
100%	85.8	60.5	67.7	87.8

Table 8: Results of the steered visual token ratio. The grey shading indicates the best results and our selected parameters.

trinsic modality imbalance evaluation. The framework aims to optimize visual instruction tuning and simplify the development process for researchers. A detailed comparison with other frameworks, such as TinyLLaVA Factory (Jia et al., 2024) and Prismatic VLMs (Karamcheti et al., 2024), is shown in Table 5. And an overview of its components is provided in Figure 7 (see Appendix A.8).

# 7 Conclusion

This introduces Modality Linear paper Representation-Steering, which significantly reduces trainable parameters while maintaining strong performance across downstream tasks by rebalancing visual and textual representations. Integrating MoReS into LLaVA models validates its effectiveness, supporting the potential of intrinsic modality rebalance for optimizing visual instruction tuning. To support future research, we present the LLaVA Steering Factory, a versatile framework enabling customizable training configurations and integrated analytical tools.

# Limitations

MoReS shows promising results, but there are areas for improvement. A more detailed analysis of its underlying mechanisms is needed to enhance interpretability and provide better insight into how it balances visual and textual representations. Additionally, further testing is required to evaluate its performance in more complex, real-world scenarios and to assess its robustness against noisy data.

#### References

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. Phi-3 technical report: A highly capable language model locally on your phone. Preprint, arXiv:2404.14219.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikoł aj Bińkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022. Flamingo: a visual language model for few-shot learning. In Advances in Neural Information Processing Systems, volume 35, pages 23716– 23736. Curran Associates, Inc.
- Afra Amini, Tiago Pimentel, Clara Meister, and Ryan Cotterell. 2023. Naturalistic causal probing for morpho-syntax. *Transactions of the Association for Computational Linguistics*, 11:384–403.
- Matan Avitan, Ryan Cotterell, Yoav Goldberg, and Shauli Ravfogel. 2024. Natural language counter-

factuals through representation surgery. *Preprint*, arXiv:2402.11355.

- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.
- Jinhe Bi, Yifan Wang, Danqi Yan, Xun Xiao, Artur Hecker, Volker Tresp, and Yunpu Ma. 2025a. Prism: Self-pruning intrinsic selection method for training-free multimodal data selection. *Preprint*, arXiv:2502.12119.
- Jinhe Bi, Danqi Yan, Yifan Wang, Wenke Huang, Haokun Chen, Guancheng Wan, Mang Ye, Xun Xiao, Hinrich Schuetze, Volker Tresp, and Yunpu Ma. 2025b. Cot-kinetics: A theoretical modeling assessing lrm reasoning process. *Preprint*, arXiv:2505.13408.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29.
- Guanzheng Chen, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. 2022. Revisiting parameterefficient tuning: Are we really there yet? In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 2612–2626, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang. 2024a. An image is worth 1/2 tokens after layer 2: Plug-andplay inference acceleration for large vision-language models. *arXiv preprint arXiv:2403.06764*.
- Yangyi Chen, Xingyao Wang, Hao Peng, and Heng Ji. 2024b. A single transformer for scalable vision-language modeling. *arXiv preprint arXiv:2407.06438*.
- Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. 2024c. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*.
- Qing Cheng, Niclas Zeller, and Daniel Cremers. 2022. Vision-based large-scale 3d semantic mapping for

autonomous driving applications. In 2022 International Conference on Robotics and Automation (ICRA), pages 9235–9242.

- Zhiqing Cui, Jiahao Yuan, Hanqing Wang, Yanshu Li, Chenxu Du, and Zhenglong Ding. 2025. Draw with thought: Unleashing multimodal reasoning for scientific diagram generation. *Preprint*, arXiv:2504.09479.
- Wenliang Dai, Nayeon Lee, Boxin Wang, Zhuoling Yang, Zihan Liu, Jon Barker, Tuomas Rintamaki, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. 2024. Nvlm: Open frontier-class multimodal llms. *Preprint*, arXiv:2409.11402.
- Haiwen Diao, Yufeng Cui, Xiaotong Li, Yueze Wang, Huchuan Lu, and Xinlong Wang. 2024. Unveiling encoder-free vision-language models. *arXiv preprint arXiv:2406.11832*.
- Guodong Du, Zitao Fang, Jing Li, Junlin Li, Runhua Jiang, Shuyang Yu, Yifei Guo, Yangneng Chen, Sim Kuan Goh, Ho-Kin Tang, Daojing He, Honghai Liu, and Min Zhang. 2025a. Neural parameter search for slimmer fine-tuned models and better transfer. *arXiv preprint arXiv:2505.18713*.
- Guodong Du, Junlin Lee, Jing Li, Runhua Jiang, Yifei Guo, Shuyang Yu, Hanting Liu, Sim Kuan Goh, Ho-Kin Tang, Daojing He, and Min Zhang. 2024. Parameter competition balancing for model merging. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS).*
- Guodong Du, Xuanning Zhou, Junlin Li, Zhuo Li, Zesheng Shi, Wanyu Lin, Ho-Kin Tang, Xiucheng Li, Fangming Liu, Wenya Wang, Min Zhang, and Jing Li. 2025b. Knowledge grafting of large language models. *arXiv preprint arXiv:2505.18502*.
- Atticus Geiger, Zhengxuan Wu, Christopher Potts, Thomas Icard, and Noah Goodman. 2024. Finding alignments between interpretable causal variables and distributed neural representations. In *Causal Learning and Reasoning*, pages 160–187. PMLR.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017a. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR).*
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017b. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913.
- Shuhao Guan and Derek Greene. 2024. Advancing post-OCR correction: A comparative study of synthetic data. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 6036–6047, Bangkok, Thailand. Association for Computational Linguistics.

- Shuhao Guan, Moule Lin, Cheng Xu, Xinyi Liu, Jinman Zhao, Jiexin Fan, Qi Xu, and Derek Greene. 2025. Prep-OCR: A complete pipeline for document image restoration and enhanced OCR accuracy. In *The 63rd Annual Meeting of the Association for Computational Linguistics*.
- Shuhao Guan, Cheng Xu, Moule Lin, and Derek Greene. 2024. Effective synthetic data and test-time adaptation for OCR correction. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 15412–15425, Miami, Florida, USA. Association for Computational Linguistics.
- Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. 2023. Hallusionbench: An advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. *arXiv preprint arXiv:2310.14566*.
- Clément Guerner, Anej Svete, Tianyu Liu, Alexander Warstadt, and Ryan Cotterell. 2023. A geometric notion of causal probing. *arXiv preprint arXiv:2307.15054*.
- Hao Guo, Zihan Ma, Zhi Zeng, Minnan Luo, Weixin Zeng, Jiuyang Tang, and Xiang Zhao. 2025. Each fake news is fake in its own way: An attribution multigranularity benchmark for multimodal fake news detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 228–236.
- Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. 2018. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3608–3617.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International conference on machine learning*, pages 2790–2799. PMLR.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Zhiyuan Hu, Yuliang Liu, Jinman Zhao, Suyuchen Wang, Yan Wang, Wei Shen, Qing Gu, Anh Tuan Luu, See-Kiong Ng, Zhiwei Jiang, and Bryan Hooi. 2024. Longrecipe: Recipe for efficient long context generalization in large language models. *Preprint*, arXiv:2409.00509.
- Wenke Huang, Jian Liang, Xianda Guo, Yiyang Fang, Guancheng Wan, Xuankun Rong, Chi Wen, Zekun Shi, Qingyun Li, Didi Zhu, et al. 2025a. Keeping yourself is important in downstream tuning

multimodal large language model. *arXiv preprint* arXiv:2503.04543.

- Wenke Huang, Jian Liang, Zekun Shi, Didi Zhu, Guancheng Wan, He Li, Bo Du, Dacheng Tao, and Mang Ye. 2025b. Learn from downstream and be yourself in multimodal large language model finetuning. In *ICML*.
- Wenke Huang, Jian Liang, Guancheng Wan, Didi Zhu, He Li, Jiawei Shao, Mang Ye, Bo Du, and Dacheng Tao. 2025c. Be confident: Uncovering overfitting in mllm multi-task tuning. In *ICML*.
- Wenke Huang, Guancheng Wan, Mang Ye, and Bo Du. 2023. Federated graph semantic and structural learning. In Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, pages 3830–3838.
- Ziwei Huang, Wanggui He, Quanyu Long, Yandi Wang, Haoyuan Li, Zhelun Yu, Fangxun Shu, Long Chan, Hao Jiang, Leilei Gan, et al. 2024. T2i-factualbench: Benchmarking the factuality of text-to-image models with knowledge-intensive concepts. *arXiv preprint arXiv:2412.04300*.
- Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709.
- Junlong Jia, Ying Hu, Xi Weng, Yiming Shi, Miao Li, Xingjian Zhang, Baichuan Zhou, Ziyu Liu, Jie Luo, Lei Huang, and Ji Wu. 2024. Tinyllava factory: A modularized codebase for small-scale large multimodal models. *Preprint*, arXiv:2405.11788.
- Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. 2024. Prismatic vlms: Investigating the design space of visually-conditioned language models. *arXiv preprint arXiv:2402.07865*.
- Karim Lasri, Tiago Pimentel, Alessandro Lenci, Thierry Poibeau, and Ryan Cotterell. 2022. Probing for the usage of grammatical number. *arXiv preprint arXiv:2204.08831*.
- Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and Victor Sanh. 2023. Obelics: An open web-scale filtered dataset of interleaved image-text documents. *Preprint*, arXiv:2306.16527.
- Bin Li, Bin Sun, Shutao Li, Encheng Chen, Hongru Liu, Yixuan Weng, Yongping Bai, and Meiling Hu. 2024a. Distinct but correct: generating diversified and entity-revised medical response. *Science China Information Sciences*, 67(3):132106.

- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. Blip-2: bootstrapping language-image pretraining with frozen image encoders and large language models. In *Proceedings of the 40th International Conference on Machine Learning*, ICML'23. JMLR.org.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2024b. Inferencetime intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36.
- Shutao Li, Bin Li, Bin Sun, and Yixuan Weng. 2024c. Towards visual-prompt temporal answer grounding in instructional video. *IEEE transactions on pattern analysis and machine intelligence*, 46(12):8836– 8853.
- Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4582– 4597, Online. Association for Computational Linguistics.
- Yanshu Li, Hongyang He, Yi Cao, Qisen Cheng, Xiang Fu, and Ruixiang Tang. 2025a. M2iv: Towards efficient and fine-grained multimodal in-context learning in large vision-language models. *Preprint*, arXiv:2504.04633.
- Yanshu Li, Tian Yun, Jianjiang Yang, Pinyuan Feng, Jinfa Huang, and Ruixiang Tang. 2025b. Taco: Enhancing multimodal in-context learning via task mapping-guided sequence configuration. *Preprint*, arXiv:2505.17098.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023b. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*.
- Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023c. Textbooks are all you need ii: phi-1.5 technical report. arXiv preprint arXiv:2309.05463.
- Yujie Lin, Ante Wang, Moye Chen, Jingyao Liu, Hao Liu, Jinsong Su, and Xinyan Xiao. 2025. Investigating inference-time scaling for chain of multimodal thought: A preliminary study. *Preprint*, arXiv:2502.11514.
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. 2022a. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. Advances in Neural Information Processing Systems, 35:1950–1965.

- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 26296–26306.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024b. Llavanext: Improved reasoning, ocr, and world knowledge.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. Visual instruction tuning. In *Advances in Neural Information Processing Systems*, volume 36, pages 34892–34916. Curran Associates, Inc.
- Tong Liu, Zhixin Lai, Gengyuan Zhang, Philip Torr, Vera Demberg, Volker Tresp, and Jindong Gu. 2024c. Multimodal pragmatic jailbreak on text-to-image models. *Preprint*, arXiv:2409.19149.
- Tong Liu, Xiao Yu, Wenxuan Zhou, Jindong Gu, and Volker Tresp. 2025. Focalpo: Enhancing preference optimizing by focusing on correct preference rankings. *Preprint*, arXiv:2501.06645.
- Wei Liu, Haozhao Wang, Jun Wang, Ruixuan Li, Xinyang Li, Yuankai Zhang, and Yang Qiu. 2023b. Mgr: Multi-generator based rationalization. *Preprint*, arXiv:2305.04492.
- Wei Liu, Haozhao Wang, Jun Wang, Ruixuan Li, Chao Yue, and YuanKai Zhang. 2022b. Fr: Folded rationalization with a unified encoder. In Advances in Neural Information Processing Systems, volume 35, pages 6954–6966. Curran Associates, Inc.
- Yilun Liu, Yunpu Ma, Shuo Chen, Zifeng Ding, Bailan He, Zhen Han, and Volker Tresp. 2024d. Perft: Parameter-efficient routed fine-tuning for mixture-of-expert model. *arXiv preprint arXiv:2411.08212*.
- Zikang Liu, Kun Zhou, Wayne Xin Zhao, Dawei Gao, Yaliang Li, and Ji-Rong Wen. 2024e. Less is more: Data value estimation for visual instruction tuning. *arXiv preprint arXiv:2403.09559*.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*.
- Pan Lu, Liang Qiu, Jiaqi Chen, Tony Xia, Yizhou Zhao, Wei Zhang, Zhou Yu, Xiaodan Liang, and Song-Chun Zhu. 2021. Iconqa: A new benchmark for abstract diagram understanding and visual language reasoning. arXiv preprint arXiv:2110.13214.
- Shilin Lu, Yanzhu Liu, and Adams Wai-Kin Kong. 2023. Tf-icon: Diffusion-based training-free cross-domain image composition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2294–2305.

- Shilin Lu, Zilan Wang, Leyang Li, Yanzhu Liu, and Adams Wai-Kin Kong. 2024a. Mace: Mass concept erasure in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6430–6440.
- Shilin Lu, Zihan Zhou, Jiayou Lu, Yuanzhi Zhu, and Adams Wai-Kin Kong. 2024b. Robust watermarking using generative priors against image editing: From benchmarking to advances. *arXiv preprint arXiv:2410.18775*.
- Zhengyang Lu and Ying Chen. 2023. Joint selfsupervised depth and optical flow estimation towards dynamic objects. *Neural Processing Letters*, 55(8):10235–10249.
- Zihan Ma, Minnan Luo, Hao Guo, Zhi Zeng, Yiran Hao, and Xiang Zhao. 2024. Event-radar: Event-driven multi-view learning for multimodal fake news detection. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5809–5821.
- Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruti Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. 2024. Mm1: Methods, analysis & insights from multimodal llm pre-training. *arXiv preprint arXiv:2403.09611*.
- Neel Nanda, Andrew Lee, and Martin Wattenberg. 2023. Emergent linear representations in world models of self-supervised sequence models. *arXiv preprint arXiv:2309.00941*.
- Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian Weller, and Bernhard Schölkopf. 2023. Controlling text-to-image diffusion by orthogonal finetuning. *Advances in Neural Information Processing Systems*, 36:79320–79362.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.
- Erica K. Shimomoto, Edison Marrese-Taylor, Hiroya Takamura, Ichiro Kobayashi, and Yusuke Miyao. 2022. A subspace-based analysis of structured and unstructured representations in image-text retrieval. In *Proceedings of the Workshop on Unimodal and Multimodal Induction of Linguistic Structures (UM-IoS)*, pages 29–44, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. 2020. Textcaps: a dataset for image captioning with reading comprehension.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards vqa models

that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326.

- Shashwat Singh, Shauli Ravfogel, Jonathan Herzig, Roee Aharoni, Ryan Cotterell, and Ponnurangam Kumaraguru. 2024. Representation surgery: Theory and practice of affine steering. *Preprint*, arXiv:2402.09631.
- Nishant Subramani, Nivedita Suresh, and Matthew E Peters. 2022. Extracting latent steering vectors from pretrained language models. *arXiv preprint arXiv:2205.05124*.
- Qwen team. 2024. Qwen2-vl.
- Guancheng Wan, Zijie Huang, Wanjia Zhao, Xiao Luo, Yizhou Sun, and Wei Wang. 2025a. Rethink graphode generalization within coupled dynamical system. In *Forty-second International Conference on Machine Learning*.
- Guancheng Wan, Zewen Liu, Xiaojun Shan, Max S.Y. Lau, B. Aditya Prakash, and Wei Jin. 2025b. Epidemiology-aware neural ode with continuous disease transmission graph. In *Forty-second International Conference on Machine Learning*.
- Guancheng Wan, Zitong Shi, Wenke Huang, Guibin Zhang, Dacheng Tao, and Mang Ye. 2025c. Energybased backdoor defense against federated graph learning. In *International Conference on Learning Representations*.
- Yaqing Wang, Sahaj Agarwal, Subhabrata Mukherjee, Xiaodong Liu, Jing Gao, Ahmed Hassan Awadallah, and Jianfeng Gao. 2022. AdaMix: Mixtureof-adaptations for parameter-efficient model tuning. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 5744–5760, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Yikun Wang, Siyin Wang, Qinyuan Cheng, Zhaoye Fei, Liang Ding, Qipeng Guo, Dacheng Tao, and Xipeng Qiu. 2025a. Visuothink: Empowering lvlm reasoning with multimodal tree search. *Preprint*, arXiv:2504.09130.
- Zhiyuan Wang, Jinhao Duan, Lu Cheng, Yue Zhang, Qingni Wang, Xiaoshuang Shi, Kaidi Xu, Hengtao Shen, and Xiaofeng Zhu. 2024a. Conu: Conformal uncertainty in large language models with correctness coverage guarantees. *Preprint*, arXiv:2407.00499.
- Zhiyuan Wang, Jinhao Duan, Chenxi Yuan, Qingyu Chen, Tianlong Chen, Yue Zhang, Ren Wang, Xiaoshuang Shi, and Kaidi Xu. 2024b. Word-sequence entropy: Towards uncertainty estimation in free-form medical question answering applications and beyond. *Preprint*, arXiv:2402.14259.
- Zhiyuan Wang, Qingni Wang, Yue Zhang, Tianlong Chen, Xiaofeng Zhu, Xiaoshuang Shi, and Kaidi Xu. 2025b. Sconu: Selective conformal uncertainty in large language models. *Preprint*, arXiv:2504.14154.

- Muling Wu, Wenhao Liu, Xiaohua Wang, Tianlong Li, Changze Lv, Zixuan Ling, Jianhao Zhu, Cenyuan Zhang, Xiaoqing Zheng, and Xuanjing Huang. 2024a. Advancing parameter efficiency in finetuning via representation editing. *arXiv preprint arXiv:2402.15179*.
- Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D. Manning, and Christopher Potts. 2024b. Reft: Representation finetuning for language models. *Preprint*, arXiv:2404.03592.
- Zhengxuan Wu, Atticus Geiger, Thomas Icard, Christopher Potts, and Noah Goodman. 2024c. Interpretability at scale: Identifying causal mechanisms in alpaca. *Advances in Neural Information Processing Systems*, 36.
- Rong Xuankun, Zhang Jianshu, He Kun, and Mang Ye. 2025. Can: Leveraging clients as navigators for generative replay in federated continual learning. In *ICML*.
- Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu, Yutong Dai, Michael S Ryoo, Shrikant Kendre, Jieyu Zhang, Can Qin, Shu Zhang, Chia-Chih Chen, Ning Yu, Juntao Tan, Tulika Manoj Awalgaonkar, Shelby Heinecke, Huan Wang, Yejin Choi, Ludwig Schmidt, Zeyuan Chen, Silvio Savarese, Juan Carlos Niebles, Caiming Xiong, and Ran Xu. 2024. xgen-mm (blip-3): A family of open large multimodal models. *Preprint*, arXiv:2408.08872.
- Ruihan Yang, Fanghua Ye, Jian Li, Siyu Yuan, Yikai Zhang, Zhaopeng Tu, Xiaolong Li, and Deqing Yang. 2025a. The lighthouse of language: Enhancing llm agents via critique-guided improvement. *Preprint*, arXiv:2503.16024.
- Shuo Yang, Siwen Luo, and Soyeon Caren Han. 2025b. Multimodal commonsense knowledge distillation for visual question answering (student abstract). In *Proceedings of the AAAI conference on artificial intelligence*, volume 39, pages 29545–29547.
- Shuo Yang, Siwen Luo, Soyeon Caren Han, and Eduard Hovy. 2025c. Magic-vqa: Multimodal and grounded inference with commonsense knowledge for visual question answering. *arXiv preprint arXiv:2503.18491*.
- Linli Yao, Lei Li, Shuhuai Ren, Lean Wang, Yuanxin Liu, Xu Sun, and Lu Hou. 2024. Deco: Decoupling token compression from semantic abstraction in multimodal large language models. *arXiv preprint arXiv:2405.20985*.
- Ting Yao, Yingwei Pan, Chong-Wah Ngo, Houqiang Li, and Tao Mei. 2015. Semi-supervised domain adaptation with subspace learning for visual recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).*

- Mang Ye, Xuankun Rong, Wenke Huang, Bo Du, Nenghai Yu, and Dacheng Tao. 2025. A survey of safety on large vision-language models: Attacks, defenses and evaluations. *arXiv preprint arXiv:2502.14881*.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. arXiv preprint arXiv:2308.02490.
- Jiahao Yuan, Dehui Du, Hao Zhang, Zixiang Di, and Usman Naseem. 2025. Reversal of thought: Enhancing large language models with preferenceguided reverse reasoning warm-up. *Preprint*, arXiv:2410.12323.
- Junrong Yue, Yifan Zhang, Chuan Qin, Bo Li, Xiaomin Lie, Xinlei Yu, Wenxin Zhang, and Zhendong Zhao. 2025. Think hierarchically, act dynamically: Hierarchical multi-modal fusion and reasoning for vision-and-language navigation. *Preprint*, arXiv:2504.16516.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.
- Zhi Zeng, Minnan Luo, Xiangzheng Kong, Huan Liu, Hao Guo, Hao Yang, Zihan Ma, and Xiang Zhao. 2024. Mitigating world biases: A multimodal multiview debiasing framework for fake news video detection. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 6492–6500.
- Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986.
- Gengyuan Zhang, Jinhe Bi, Jindong Gu, Yanyu Chen, and Volker Tresp. 2023a. Spot! revisiting videolanguage models for event understanding. *Preprint*, arXiv:2311.12919.
- Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024a. Tinyllama: An open-source small language model. *Preprint*, arXiv:2401.02385.
- Xueyan Zhang, Jinman Zhao, Zhifei Yang, Yibo Zhong, Shuhao Guan, Linbo Cao, and Yining Wang. 2025. Uora: Uniform orthogonal reinitialization adaptation in parameter-efficient fine-tuning of large models. *Preprint*, arXiv:2505.20154.
- Yao Zhang, Zijian Ma, Yunpu Ma, Zhen Han, Yu Wu, and Volker Tresp. 2024b. Webpilot: A versatile and autonomous multi-agent system for web task execution with strategic exploration. *Preprint*, arXiv:2408.15978.

- Yi-Kai Zhang, Shiyin Lu, Yang Li, Yanqing Ma, Qing-Guo Chen, Zhao Xu, Weihua Luo, Kaifu Zhang, De-Chuan Zhan, and Han-Jia Ye. 2024c. Wings: Learning multimodal llms without text-only forgetting. arXiv preprint arXiv:2406.03496.
- Zheyu Zhang, Han Yang, Bolei Ma, David Rügamer, and Ercong Nie. 2023b. Baby's cothought: Leveraging large language models for enhanced reasoning in compact models. *Preprint*, arXiv:2308.01684.
- Jinman Zhao and Xueyan Zhang. 2024. Large language model is not a (multilingual) compositional relation reasoner. In *First Conference on Language Modeling*.
- Xuanle Zhao, Xianzhen Luo, Qi Shi, Chi Chen, Shuo Wang, Wanxiang Che, Zhiyuan Liu, and Maosong Sun. 2025. Chartcoder: Advancing multimodal large language model for chart-to-code generation. *arXiv preprint arXiv:2501.06598*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. 2023a. Judging Ilm-as-a-judge with mt-bench and chatbot arena. In Advances in Neural Information Processing Systems, volume 36, pages 46595–46623. Curran Associates, Inc.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023b. Judging llm-as-a-judge with mt-bench and chatbot arena. *Preprint*, arXiv:2306.05685.
- Lin Zhong, Jun Zeng, Ziwei Wang, Wei Zhou, and Junhao Wen. 2024. Scfl: Spatio-temporal consistency federated learning for next poi recommendation. *Information Processing & Management*, 61(6):103852.
- Lin Zhong, Jun Zeng, Yang Yu, Hongjin Tao, Wenying Jiang, and Luxi Cheng. 2023. A text matching model based on dynamic multi-mask and augmented adversarial. *Expert Systems*, 40(2):e13165.
- Baichuan Zhou, Ying Hu, Xi Weng, Junlong Jia, Jie Luo, Xien Liu, Ji Wu, and Lei Huang. 2024a. Tinyllava: A framework of small-scale large multimodal models. *Preprint*, arXiv:2402.14289.
- Xingcheng Zhou, Xuyuan Han, Feng Yang, Yunpu Ma, and Alois C. Knoll. 2025. Opendrivevla: Towards end-to-end autonomous driving with large vision language action model. *Preprint*, arXiv:2503.23463.
- Xiongtao Zhou, Jie He, Yuhua Ke, Guangyao Zhu, Víctor Gutiérrez-Basulto, and Jeff Z Pan. 2024b. An empirical study on parameter-efficient fine-tuning for multimodal large language models. *arXiv preprint arXiv:2406.05130*.

- zhou changshi, Feng Luan, hujiarui, Shaoqiang Meng, Zhipeng Wang, Yanchao Dong, Yanmin Zhou, and Bin He. 2025. Learning efficient robotic garment manipulation with standardization. In *ICML*.
- Xingyu Zhu, Beier Zhu, Yi Tan, Shuo Wang, Yanbin Hao, and Hanwang Zhang. 2024. Selective visionlanguage subspace projection for few-shot clip. *arXiv* preprint arXiv:2407.16977.
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. 2023. Representation engineering: A topdown approach to ai transparency. *arXiv preprint arXiv:2310.01405*.

# Acknowledgements

The authors gratefully acknowledge the scientific support and HPC resources provided by the Erlangen National High Performance Computing Center (NHR@FAU) of the Friedrich-Alexander-Universität Erlangen-Nürnberg (FAU) under the NHR projects b207dd and b211dd. NHR funding is provided by federal and Bavarian state authorities. The NHR@FAU hardware is partially funded by the German Research Foundation (DFG) – 440719683.

# **A** Appendix

# A.1 Experiment Settings

## A.1.1 LLaVA Steering Architectures

As illustrated in Figure 3, the architecture of the LLaVA Steering models (3B/7B/13B) consists of three essential components: a vision encoder, a vision connector responsible for projecting visual representations into a shared latent space, and a multi-scale LLM. The three modules are introduced below.

In our experiments, we utilize the Phi-2 2.7B model (Li et al., 2023c) alongside Vicuna v1.5 (7B and 13B) (Zheng et al., 2023b), sourced from our factory, to evaluate the generalizability of our approach across models of varying scales. For vision encoding, we employ CLIP ViT-L/14 336px (Radford et al., 2021) and SigLIP-SO400M-Patch14-384 (Zhai et al., 2023), while a two-layer MLP serves as the connector. Given the inefficiencies of Qformer in training and its tendency to introduce cumulative deficiencies in visual semantics (Yao et al., 2024), it has been largely replaced by more advanced architectures, such as the BLIP series (Xue et al., 2024), Qwen-VL series (team, 2024), and InternVL series (Chen et al., 2024c), which were previously reliant on Qformer.

#### A.1.2 Baseline Training Methods

For comparison, four widely adopted PEFT methods (Adapter, LoRA, OFT and IA3) are selected as baselines. These methods establish a comparative framework to assess both the performance and efficiency of our proposed approach. Essentially, our MoReS method replaces these four PEFT methods during visual instruction tuning in LLaVA Steering. Adapter: Building on the framework of efficient fine-tuning (Houlsby et al., 2019), we introduce adapter layers within Transformer blocks. These layers consist of a down-projection matrix  $\mathbf{W}_{\text{down}} \in \mathbb{R}^{r \times d}$ , a non-linear activation function  $\sigma(\cdot)$ , and an up-projection matrix  $\mathbf{W}_{up} \in \mathbb{R}^{d \times r}$ , where d is the hidden layer dimension and r is the bottleneck dimension. The adapter output is computed as:

$$Adapter(\mathbf{x}) = \mathbf{W}_{up}\sigma(\mathbf{W}_{down}\mathbf{x}) + \mathbf{x}, \quad (5)$$

where the residual connection (+x) preserves the pre-trained model's knowledge. This formulation enables efficient parameter updates during fine-tuning, offering a balance between computational efficiency and adaptation capacity while minimally increasing the model's complexity.

**LoRA:** We employ the low-rank adaptation method (LoRA) proposed by (Hu et al., 2021), which efficiently updates the network's weights with a minimal parameter footprint by leveraging a low-rank decomposition strategy. For a pre-trained weight matrix  $W_0 \in \mathbb{R}^{d \times k}$ , the weight update is achieved through the addition of a low-rank decomposition, as shown in Equation 6:

$$W_0 + \Delta W = W_0 + BA \tag{6}$$

where  $B \in \mathbb{R}^{d \times r}$  and  $A \in \mathbb{R}^{r \times k}$  are trainable low-rank matrices, and  $r \ll \min(d, k)$ .

**OFT:** We utilize the Orthogonal Finetuning (OFT) method, which efficiently fine-tunes pre-trained models by optimizing a constrained orthogonal transformation matrix (Qiu et al., 2023). For a pre-trained weight matrix  $W_0 \in \mathbb{R}^{d \times n}$ , OFT modifies the forward pass by introducing an orthogonal matrix  $R \in \mathbb{R}^{d \times d}$ , as illustrated in Equation 7:

$$z = W^{\top} x = (R \cdot W_0)^{\top} x \tag{7}$$

where R is initialized as an identity matrix I to ensure that fine-tuning starts from the pre-trained weights.

**IA3:** Building on the framework established by (Liu et al., 2022a), we introduce three vectors  $v_k \in$ 

 $\mathbb{R}^{d_k}$ ,  $v_v \in \mathbb{R}^{d_v}$ , and  $v_{ff} \in \mathbb{R}^{d_{ff}}$  into the attention mechanism. The attention output is computed as:

Attention = softmax 
$$\left(\frac{Q(v_k \odot K^T)}{\sqrt{d_k}}\right) (v_v \odot V),$$
(8)

where  $\odot$  denotes multiplication by element.

#### A.2 Benchmarks Overview

Recent advances in deep learning (Wan et al., 2025a,b,c; Huang et al., 2023; Lu et al., 2024b,a, 2023; Lu and Chen, 2023; zhou changshi et al., 2025; Liu et al., 2023b, 2022b; Guan et al., 2024; Guan and Greene, 2024; Guan et al., 2025) have led to the emergence of large language models (LLMs) (Zhang et al., 2024b; Liu et al., 2024d; Ye et al., 2025; Xuankun et al., 2025; Yang et al., 2025a; Zhang et al., 2023b; Bi et al., 2025b; Yuan et al., 2025; Wang et al., 2024a,b, 2025b; Lin et al., 2025; Zhao and Zhang, 2024; Hu et al., 2024; Du et al., 2024, 2025b,a; Zhong et al., 2023, 2024), which demonstrate remarkable capabilities across a broad range of NLP tasks. Building upon these successes, multimodal large language models (MLLMs) have been developed to extend this capability to visuallinguistic tasks by integrating image understanding with natural language reasoning. To comprehensively evaluate the performance of such models, various benchmarks have been proposed that assess different aspects of multimodal understanding (Zhou et al., 2025; Cheng et al., 2022; Liu et al., 2024c, 2025; Cui et al., 2025; Yang et al., 2025b,c; Li et al., 2025a,b; Huang et al., 2024; Zhang et al., 2023a; Ma et al., 2024; Zeng et al., 2024; Guo et al., 2025), including perception, reasoning, grounding, and hallucination. In parallel, several recent works have emphasized the need for efficient multimodal adaptation and task-specific evaluation, highlighting the importance of standardized benchmarks. Below, we provide a brief overview of the benchmarks used in our study.

**VQAv2** (Goyal et al., 2017b): A benchmark for evaluating visual perception through open-ended short answers to visual questions.

**GQA** (Hudson and Manning, 2019): A dataset for assessing visual reasoning and question answering. **VizWiz** (Gurari et al., 2018): Consists of 8,000 images designed for zero-shot generalization in visual questions posed by visually impaired individuals.

ScienceQA (Lu et al., 2022): A benchmark

focusing on zero-shot scientific question answering with multiple-choice questions.

**TextVQA** (Singh et al., 2019): Evaluates performance on text-rich visual questions.

**MM-Vet** (Yu et al., 2023): Assesses the model's ability to engage in visual conversations, with correctness and helpfulness evaluated by GPT-4.

**POPE** (Li et al., 2023b): Quantifies hallucination in MLLMs.

**MMMU** (Yue et al., 2024): Evaluates core multimodal skills, including perception, knowledge, and reasoning.

#### A.3 Hallucination Evaluation Details

POPE (Li et al., 2023b) specifically focuses on object hallucination, using accuracy (*Acc*) as the primary evaluation metric. By assessing whether the generated outputs accurately correspond to objects present in the visual input, POPE provides a clear measure of hallucination mitigation.

HallucinationBench (Guan et al., 2023) offers a broader assessment by covering diverse topics and visual modalities. This benchmark includes two categories of questions: (1) *Visual Dependent (VD) Questions*, which require detailed understanding of the visual input for correct responses, and (2) *Visual Supplement (VS) Questions*, where answers depend on contextual visual support rather than direct visual grounding.

To evaluate model performance comprehensively, we focus on three main metrics: *Hard Acc*, which assesses correctness based on strict adherence to the visual context; *Figure Acc*, measuring accuracy on a per-figure basis; and *Question Acc*, evaluating the overall accuracy across all questions.

#### A.4 Theoretical Justification

Let  $x_{\text{text}} \in \mathbb{R}^{d_t}$  be the text input embedding,  $x_{\text{image}} \in \mathbb{R}^{d_v}$  be the visual input embedding,  $R_{\text{text}} \in \mathbb{R}^D$  be the hidden representation for text, and  $R_{\text{image}} \in \mathbb{R}^D$  be the hidden representation for the visual input. Define  $W_q, W_k, W_v \in \mathbb{R}^{D \times D}$  as the query, key, and value projection matrices, and  $W_o \in \mathbb{R}^{D \times D}$  as the output projection matrix. Let  $A \in \mathbb{R}^{N \times N}$  represent the attention matrix, and  $y \in \mathbb{R}^V$  be the output logits.

We present a theoretical analysis of the MoReS transformation and its effect on attention redistribution in multimodal models. The hidden representations for text and image inputs are computed as:

$$h_{\text{text}} = f_{\text{text}}(x_{\text{text}}), \quad h_{\text{image}} = f_{\text{image}}(x_{\text{image}})$$
 (9)

where  $f_{\text{text}}$  and  $f_{\text{image}}$  are encoding functions. The attention mechanism is characterized by scores:

$$A_{ij} = \operatorname{softmax}\left(\frac{(h_i W_q)(h_j W_k)^T}{\sqrt{D}}\right) \qquad (10)$$

with  $W_q, W_k \in \mathbb{R}^{D \times D}$  being query and key projection matrices. Output generation follows:

$$y = W_o(C_{\text{text}} + C_{\text{image}}) \tag{11}$$

where  $C_{\text{text}} = \sum_{i} A_{i,\text{text}}(h_i W_v)$  and  $C_{\text{image}} = \sum_{i} A_{i,\text{image}}(h_i W_v)$ .

The core of our approach is the MoReS transformation, defined as:

$$MoReS(h) = W_{up} \cdot \phi(h), \qquad (12)$$

where 
$$\phi(h) = \text{Linear}(h) - W_{\text{down}}h$$
 (13)

Here,  $W_{up} \in \mathbb{R}^{D \times d}$ ,  $W_{down} \in \mathbb{R}^{d \times D}$ , and d < D. When applied to the image representation, we obtain  $h'_{image} = MoReS(h_{image}) + h_{image}$ , leading to updated attention scores:

$$A_{i,\text{image}}' = \text{softmax}\left(\frac{(h_i W_q)(h_{\text{image}}' W_k)^T}{\sqrt{D}}\right)$$
(14)

This transformation is key to redistributing attention towards visual inputs. The effect of MoReS on the output can be quantified by examining the change magnitude:

$$\|\Delta y\|_{2} = \|W_{o}(C'_{\text{image}} - C_{\text{image}})\|_{2}$$
(15)

$$\leq \|W_o\|_2 \|C'_{\text{image}} - C_{\text{image}}\|_2 \qquad (16)$$

where  $C'_{\text{image}} = \sum_i A'_{i,\text{image}}(h'_{\text{image}}W_v)$ . The significance of this change stems from the MoReS transformation's ability to amplify key visual features. Specifically,  $\phi(h)$  extracts salient visual information in a subspace, which is then amplified by  $W_{\text{up}}$  in the original space. This process ensures  $\|h'_{\text{image}}\|_2 > \|h_{\text{image}}\|_2$ , leading to increased  $A'_{i,\text{image}}$  values for relevant visual features and larger magnitudes for  $(h'_{\text{image}}W_v)$  terms in  $C'_{\text{image}}$ . To ensure stability while allowing for this significant attention redistribution, we consider the Lipschitz continuity of the model:

$$\|f(h'_{\text{image}}) - f(h_{\text{image}})\|_2 \le L \|h'_{\text{image}} - h_{\text{image}}\|_2$$
(17)

where L is the Lipschitz constant. This property bounds the change in the model's output, guaranteeing that the attention redistribution, while substantial, remains controlled and does not destabilize the overall model behavior.

A key advantage of the MoReS approach lies in its parameter efficiency. The transformation introduces O(Dd) parameters, primarily from  $W_{up}$ ,  $W_{down}$ , and the linear transformation in  $\phi(h)$ . This is significantly less than the  $O(D^2)$  parameters required for fine-tuning all attention matrices in traditional approaches. The reduction in trainable parameters not only makes the optimization process more efficient but also mitigates the risk of overfitting, especially in scenarios with limited training data.

In conclusion, our theoretical analysis demonstrates that our MoReS effectively redistributes attention to visual inputs by operating in a carefully chosen subspace. This approach achieves a significant change in output generation while maintaining model stability and requiring fewer parameters than full fine-tuning, offering a balance between effectiveness and efficiency in enhancing visual understanding in MLLMs.

### A.5 Implementation Detail



Figure 5: MoReS module flowchart.

Regarding the implementation, we have adopted a highly modular design for the LLM, integrating it with MoReS to enable precise steering at specific

Steering Ratio	Dataset	LR	Accuracy (%)
0-1%	SQA	2e-4	86.7
1-25%	SQA	2e-4	87.7
25-50%	SQA	2e-4	84.1
50-100%	SQA	2e-4	81.0
0–1%	VizWiz	2e-4	58.1
1–25%	VizWiz	2e-4	58.6
25-50%	VizWiz	2e-4	58.3
50-100%	VizWiz	2e-4	59.1
0-1%	IconQA-Text	2e-4	81.2
1–25%	IconQA-Text	2e-4	81.3
25-50%	IconQA-Text	2e-4	76.4
50-100%	IconQA-Text	2e-4	77.1
0–1%	IconQA-Blank	2e-4	93.9
1–25%	IconQA-Blank	2e-4	91.2
25-50%	IconQA-Blank	2e-4	87.4
50-100%	IconQA-Blank	2e-4	84.1

Table 9: Effect of varying steering ratio across different tasks.

token locations. This modular approach ensures that the steering process operates with minimal computational overhead, making it both efficient and scalable. Additionally, the modular nature of this design allows for seamless integration with existing architectures and enables easy customization of steering strategies tailored to specific downstream tasks. To provide further clarity, we include a MoReS module flowchart (Figure 5) and an UML diagram (Figure 6) here, which detail the implementation process.



Figure 6: The UML diagram for MoReS

#### A.6 Full Attention Maps

In this section, we provide the attention maps (Figure 8) during the decoding process across each layer. Notably, the distribution of visual attention remains sparse in these layers, with only a few tokens carrying the majority of the attention. This sparsity presents an opportunity for token pruning strategies, which can be leveraged to reduce inference overhead and improve computational efficiency. By selectively pruning tokens with lower attention scores, unnecessary computations can be avoided, leading to faster and more efficient inference while maintaining the essential information needed for accurate predictions.

#### A.7 Runtime Overhead

Unlike LoRA, where the learned weights can be merged into the model's original parameters to achieve zero computational overhead during inference, MoReS requires the linear transformation layers to remain in the computation graph of the MLLM. While this introduces a small overhead, we have worked to minimize it effectively.

To mitigate runtime overhead, we performed several experiments focusing on key factors: Subspace Rank, Steered Visual Token Rate and Steering Layer Configuration. These experiments helped us reduce the additional computational burden. Specifically, by choosing a 1% Steered Visual Token Rate, a Subspace Rank of 1, and employing a sparse Steering Layer Configuration, we achieved the minimum runtime overhead of about 0.08 seconds each sample. This is significantly lower compared to other PEFT methods, such as Adapter (0.3 seconds) and OFT (2.8 seconds).

#### A.8 LLaVA Steering Factory

An overview of the main components of the LLaVA Steering Factory is provided in Figure 7.



Figure 7: Architectural overview of the proposed LLaVA Steering Factory: A Modular Codebase for MLLMs.

# A.9 Impact of Removing Linear Transformations

As shown in Table 10 and 11, we conducted experiments applying MoReS with different fixed intervals and also evaluated its performance when applied exclusively to the shallow, middle, and deep layers. These experiments highlight that the choice of steering layers can effectively balance computational efficiency and performance. We suggest that, when using MoReS, it is optimal to apply it to all layers initially to achieve the best performance. Then, by skipping fixed intervals, we can further reduce inference overhead while maintaining performance. Regarding the choice of shallow, middle, and deep layers, we found that applying MoReS to the deep layers yields better performance. We believe that deep layers encode more abstract concepts and are more suitable for steering in the subspace.

Steering Layer	VQAv2	GQA	TextVQA	SciQA-IMG	POPE	MM-Vet	MMMU	Avg
[0,2,4,]	74.1	52.0	48.3	71.6	87.1	32.8	35.3	57.3
[0,3,6,]	74.1	51.7	48.1	70.7	87.0	32.7	33.2	56.8
[0,4,8,]	74.1	51.9	48.5	71.2	87.2	31.5	34.4	57.0
All Layer	74.0	51.6	49.3	71.6	87.2	33.3	34.4	57.3

Table 10: Performance of different steering layer strategies across benchmarks.

Steering Layer	VQAv2	GQA	TextVQA	SciQA-IMG	POPE	MM-Vet	MMMU	Avg
Shallow (0-15)	74.3	51.6	48.6	70.3	87.5	34.9	34.4	57.3
Middle (8-23)	74.3	52.3	48.3	71.5	87.1	32.0	32.6	56.9
Deep (16-31)	74.2	51.5	48.2	71.8	87.1	33.3	36.7	57.7
All Layer	74.0	51.6	49.3	71.6	87.2	33.3	34.4	57.3

Table 11: Performance comparison of shallow, middle, and deep steering layers.



Figure 8: Full Attention Maps of Each Layer