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

**Anonymous ACL submission** 

#### Abstract

001 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 016 these modalities can significantly reduce trainable parameters, inspiring further optimiza-017 tion 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 vi-022 sual 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 027 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-the-art models and evaluating intrinsic modality imbalance. This open-source project facilitates a deeper understanding of MLLMs within the research 037 community.

### 1 Introduction

039

042

Recent advancements in Multimodal Large Language Models (MLLMs) (Liu et al., 2024b; Xue et al., 2024; Zhou et al., 2024a; Chen et al., 2023) have demonstrated impressive capabilities across a variety of visual downstream tasks. 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). 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

083

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., 2023).

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., 2024c; 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, such as LoRA (Hu et al., 2021), aiming to reduce the number 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 finetuning.

This raises a critical question: is there any further possibility to reduce more trainable parameters so that the visual instruction tuning can be further im-



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.

proved? 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., 2022), MoReS focuses solely on steering the visual representations. Specifically, our approach freezes the entire LLM during visual instruction tuning to preserve 113 its capabilities in the textual modality. Instead of 114 fine-tuning the full model, we introduce a simple 115 linear transformation to steer visual representations 116

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.

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

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, multi-150 scale LLMs, and diverse connectors, while offering 151 customizable training configurations across a vari-152 ety of downstream tasks. This framework simpli-153 fies pretraining and visual instruction tuning, min-154 imizing the coding effort. Additionally, we have 155 integrated our attention score distribution analy-156 sis into the LLaVA Steering Factory, providing a 157 valuable tool to the research community for further 158 studying intrinsic modality imbalance in MLLMs. 159 Our work makes the following key contributions to the field of MLLMs: 161

162

163

165

166

170

171

172

173

174

175

176

177

178

179

181

182

184

188

190

191

- 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.
- 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.
- 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., 2023)

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. 199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

**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., 2024; 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 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 representationtuning 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 repre-



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.

sentations, thereby demonstrating superior parameter efficiency compared to baseline PEFT methods (Hu et al., 2021; Houlsby et al., 2019; Liu et al., 2022; Qiu et al., 2023).

#### **3** Intrinsic Modality Imbalance

253

258

259

260

261

264

265

268

271

272

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 representa-

tion,  $R_{\text{sys}}$  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.

In our experiments comparing various existing PEFT methods and full fine-tuning, IA3 (Liu et al., 2022) 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



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.

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

326

327

328

336

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 338 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 342 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). 346 Building upon the intervention mechanisms de-347 scribed in Geiger et al. (2024) and Guerner et al. (2023), we introduce a simple linear transformation that steers visual representations within sub-350

space 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 the fidelity of the representation and

375

351

352

353

355



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

ensure a bijective mapping between spaces, we impose the following constraint  $W_{\text{down}}W_{\text{up}}^T = I_D$ . 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 5.6.

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 realworld 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 composed models across several popular and prestigious datasets.

#### 5 **Experiments**

376

379

381

394

400

401 402

403

404

We incorporate MoReS into each layer of the LLM 398 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 Ap-405 pendix A.1. 406

#### 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.

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

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 instruciton tuning stage.

#### Multi-scale Data Fine-tuning 5.3

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

Model	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	188.74M	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
LLavA Steering	IA3	0.492M	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. For the TP\* metric in this evaluation, we focus solely on the trainable parameters within the LLM. While different training strategies are applied to the vision encoder and connector across various recipes, we maintain a consistent training recipe for all models and benchmarks to ensure comparability

Model	Method	TP*	SciQA-IMG	VizWiz	IconQA-txt	IconQA-blank
	Adapter	83M	92.3	62.9	93.5	95.8
	LoRA	188.7M	93.9	61.6	93.9	96.5
LLaVA Steering-3B	OFT	39.32M	86.3	42.0	87.8	42.0
	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
	LoRA	319.8M	87.6	60.6	77.7	70.2
LLaVA Steering-7B	OFT	100.7M	78.3	55.1	19.4	22.7
	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
	LoRA	500.7M	92.1	62.0	80.2	73.2
LLaVA Steering-13B	OFT	196.6M	82.7	59.5	3.4	22.3
	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: Results of Task-Specific Fine-tuning, where higher values correspond to better performance.

sample 1K, 5K, and 10K data points from each dataset, defining these as small-scale, medium-scale, 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 smallscale tasks, MoReS outperforms full fine-tuning performance while using only 575 times fewer parameters than LoRA and 8,475 fewer than full finetuning. 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

456

457

458

459

460

461

462

463

464

465

466

467

468

470

471

472

473

474 MoReS preserves 100% of the pre-trained world
475 knowledge in the LLM by neither modifying its
476 parameters nor interfering with textual token in477 ference. This design allows MoReS to excel in
478 understanding both visual and textual information.
479 Unlike many existing methods, which often al-

Scale	Method	TP*	SciQA-IMG	VizWiz	IconQA
	FT	2.78B	33.8	51.2	68.1
	Adapter	83M	81.0	57.4	72.4
C	LoRA	188.74M	84.0	58.5	74.2
Small	OFT	39.32M	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
	FT	2.78B	78.2	58.9	92.2
	Adapter	83M	92.1	60.6	93.2
N . 11	LoRA	188.74M	92.9	60.5	92.7
Medium	OFT	39.32M	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
	FT	2.78B	88.9	59.4	95.7
	Adapter	83M	92.4	61.3	95.2
Large	LoRA	188.74M	93.9	61.8	96.0
	OFT	39.32M	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

ter model weights and risk degrading pre-trained knowledge (Zhang et al., 2024b), MoReS employs a representation-steering approach to selectively enhance the performance of the visual modality.

Text-only 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.

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. 484

485

486

487

488

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

Table 5: Comparison of functionality across different factories.

#### 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.

	Metric	Full	LoRA	Adapter	OFT	IA3	MoReS
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↑	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.

#### 5.6 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 512 found that a rank of 1 achieves the highest average 513 performance of 81.8 across four visual tasks while also requiring the fewest parameters, specifically 515 0.164M. Second, regarding the steered visual token 516 ratio, we varied this parameter from 100% (dense 517 steering) to 1% (sparse steering). The results indicate that a ratio of 1% is optimal, yielding the best or near-optimal performance on four bench-520 marks while also significantly reducing inference 521 overhead due to its sparse steering approach. 522

6 LLaVA Steering Factory

524 The LLaVA Steering Factory addresses the need for525 a comprehensive framework to systematically ana-

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.

Steered 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.

lyze 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 intrinsic 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 paper Linear 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.

552

553

554

555

498

499

500

504

490

- 505
- 506 507 508

509

510

566

567

568

569

570

571

572

573

574 575

576

577

578

579

580

581

582

583

584

585

586

588

589

592

593 594

595

596

597

599

606

607

610

611

612

557 MoReS shows promising results, but there are ar-558 eas for improvement. A more detailed analysis of 559 its underlying mechanisms is needed to enhance 560 interpretability and provide better insight into how 561 it balances visual and textual representations. Ad-562 ditionally, further testing is required to evaluate 563 its performance in more complex, real-world sce-564 narios and to assess its robustness against noisy 565 data.

#### References

Limitations

- 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. 613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

- 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 counterfactuals 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*.
- 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*.
- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023. Sharegpt4v: Improving large multimodal models with better captions. *arXiv preprint arXiv:2311.12793*.
- Yangyi Chen, Xingyao Wang, Hao Peng, and Heng Ji. 2024b. A single transformer for scalable vision-language modeling. *arXiv preprint arXiv:2407.06438.*

- 674 687 688
- 6999 7000 7011 7022 7033 7044 705 7066 7077 7088 7099 7110 7111 7122 7133 7144
- 714 715 716 717 718

72

722 723

- 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*.
- 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*.
- 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.
- 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*.
- 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*. 724

725

726

727

728

729

730

732

733

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

758

759

760

761

763

764

765

766

767

768

769

770

771

772

774

776

778

- 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.
- 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. 2024. Inferencetime intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36.
- 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.
- 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*.

- 780 781
- 790

- 798
- 807
- 810 811 812
- 813 815

816

817 818

- 819 820 821
- 822

823

824 825

826

827

829

832

833

- Haokun Liu, Derek Tam, Mohammed Mugeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A Raffel. 2022. 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. 2023. Visual instruction tuning. In Advances in Neural Information Processing Systems, volume 36, pages 34892-34916. Curran Associates, Inc.
- Zikang Liu, Kun Zhou, Wayne Xin Zhao, Dawei Gao, Yaliang Li, and Ji-Rong Wen. 2024c. 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.
- 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. 834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

- 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 vga 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.

- 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.
- Muling Wu, Wenhao Liu, Xiaohua Wang, Tianlong Li, Changze Lv, Zixuan Ling, Jianhao Zhu, Cenyuan Zhang, Xiaoqing Zheng, and Xuanjing Huang. Advancing parameter efficiency in fine-2024a. tuning 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.
- 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

951

- 952 953 954 955 956 957
- 958 959 960
- 961 962

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.

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*.

895

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916 917

918

919

920

921

925

929 930

931

932

933 934

935

936 937

940

941

- 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).*
- 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.
- 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.
- 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.
- Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024a. Tinyllama: An open-source small language model. *Preprint*, arXiv:2401.02385.
- 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. 2024b. Wings: Learning multimodal llms without text-only forget-ting. *arXiv preprint arXiv:2406.03496*.
- 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.
- 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.

- 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*.
- 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*.

# 964

967

969

970

971

974

975

976

977

978

981

987

991

995

998

1001

1002

1003

1004

1006

1007

1008

1010

# Appendix A.1 Experiment Settings

А

### 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 finetuning, 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., 2022), 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

VQAv2 (Goyal et al., 2017b): A benchmark for 1042 evaluating visual perception through open-ended 1043 short answers to visual questions. GOA (Hudson 1044 and Manning, 2019): A dataset for assessing vi-1045 sual reasoning and question answering. VizWiz 1046 (Gurari et al., 2018): Consists of 8,000 images de-1047 signed for zero-shot generalization in visual ques-1048 tions posed by visually impaired individuals. Sci-1049 enceQA (Lu et al., 2022): A benchmark focus-1050 ing on zero-shot scientific question answering with 1051 multiple-choice questions. TextVQA (Singh et al., 1052 2019): Evaluates performance on text-rich visual questions. MM-Vet (Yu et al., 2023): Assesses 1054

13

1039 1040

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1036

1037

1038

1055the model's ability to engage in visual conversa-1056tions, with correctness and helpfulness evaluated1057by GPT-4. **POPE** (Li et al., 2023b): Quantifies hal-1058lucination in MLLMs. **MMMU** (Yue et al., 2024):1059Evaluates core multimodal skills, including percep-1060tion, knowledge, and reasoning.

#### A.3 Hallucination Evaluation Details

1061

1063

1064

1067

1068

1069

1071

1073

1074

1075

1076

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1093

1094

1095

1096

1098

1099

1100

1101

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:

$$A_{ij} = \operatorname{softmax}\left(\frac{(h_i W_q)(h_j W_k)^T}{\sqrt{D}}\right) \quad (10) \quad 110$$

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

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

1103

1104

1108

1109

1

1116

1117

1118

1119

1120

1122

1132

1133

1134

where 
$$C_{\text{text}} = \sum_{i} A_{i,\text{text}}(h_i W_v)$$
 and  $C_{\text{image}} = 1106$   
 $\sum_{i} A_{i,\text{image}}(h_i W_v)$ . 1107

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

$$MoReS(h) = W_{up} \cdot \phi(h), \qquad (12) \qquad 111$$

where 
$$\phi(h) = \text{Linear}(h) - W_{\text{down}}h$$
 (13) 111

Here,  $W_{up} \in \mathbb{R}^{D \times d}$ ,  $W_{down} \in \mathbb{R}^{d \times D}$ , and d < 1112 D. When applied to the image representation, we obtain  $h'_{image} = MoReS(h_{image}) + h_{image}$ , leading 1114 to updated attention scores: 1115

$$A_{i,\text{image}}' = \text{softmax}\left(\frac{(h_i W_q)(h_{\text{image}}' W_k)^T}{\sqrt{D}}\right)$$
(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) 1

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

where  $C'_{\text{image}} = \sum_i A'_{i,\text{image}}(h'_{\text{image}}W_v)$ . The 1123 significance of this change stems from the MoReS 1124 transformation's ability to amplify key visual fea-1125 tures. Specifically,  $\phi(h)$  extracts salient visual in-1126 formation in a subspace, which is then amplified 1127 by  $W_{up}$  in the original space. This process en-1128 sures  $\|h'_{\text{image}}\|_2 > \|h_{\text{image}}\|_2$ , leading to increased 1129  $A'_{i,\text{image}}$  values for relevant visual features and 1130 larger magnitudes for  $(h'_{\text{image}}W_v)$  terms in  $C'_{\text{image}}$ . 1131

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) 1135

1174

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1207

1175

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.

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

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 1162 a highly modular design for the LLM, integrating 1163 it with MoReS to enable precise steering at specific 1164 token locations. This modular approach ensures 1165 that the steering process operates with minimal 1166 1167 computational overhead, making it both efficient and scalable. Additionally, the modular nature of 1168 this design allows for seamless integration with 1169 existing architectures and enables easy customiza-1170 tion of steering strategies tailored to specific down-1171

stream 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).
- 1210 A.8 LLaVA Steering Factory

1211An overview of the main components of the LLaVA1212Steering 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

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

As shown in Table 9 and 10, 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 9: 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 10: Performance comparison of shallow, middle, and deep steering layers.



Figure 8: Full Attention Maps of Each Layer