



VISUAL INSTRUCTION TUNING WITH 500X FEWER PARAMETERS THROUGH MODALITY LINEAR REPRESENTATION-STEERING

Anonymous authors

Paper under double-blind review

ABSTRACT

Multimodal Large Language Models (MLLMs) have significantly advanced visual tasks by integrating visual representations into large language models (LLMs). The textual modality, inherited from LLMs, equips MLLMs with abilities like instruction following and in-context learning. In contrast, the visual modality enhances performance in downstream tasks by leveraging rich semantic content, spatial information, and grounding capabilities. These intrinsic modalities work synergistically across various visual tasks. Our research initially reveals a persistent imbalance between these modalities, with text often dominating output generation during visual instruction tuning. This imbalance occurs when using both full fine-tuning and parameter-efficient fine-tuning (PEFT) methods. We then found that re-balancing these modalities can significantly reduce the number of trainable parameters required, inspiring a direction for further optimizing visual instruction tuning. Hence, in this paper, we introduce Modality Linear Representation-Steering (MoReS) to achieve the goal. MoReS effectively re-balances the intrinsic modalities throughout the model, where the key idea is to steer visual representations through linear transformations in the visual subspace across each model layer. To validate our solution, we composed LLaVA Steering, a suite of models integrated with the proposed MoReS method. Evaluation results show that the composed LLaVA Steering models require, on average, 500 times fewer trainable parameters than LoRA needs while still achieving comparable performance across three visual benchmarks and eight visual question-answering tasks. Last, we present the LLaVA Steering Factory, an in-house developed platform that enables researchers to quickly customize various MLLMs with component-based architecture for seamlessly integrating state-of-the-art models, and evaluate their intrinsic modality imbalance. This open-source project enriches the research community to gain a deeper understanding of MLLMs.

1 INTRODUCTION

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., 2024; 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

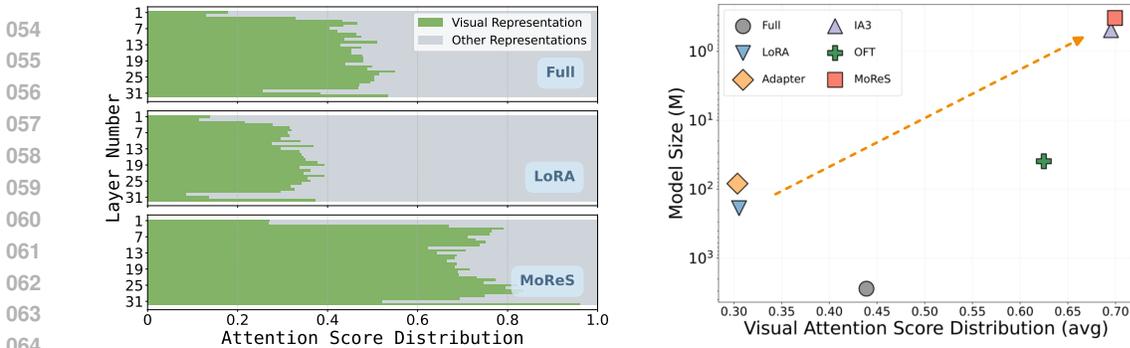


Figure 1: **Left:** Attention score distributions across layers for three MLLM fine-tuning methods (Full, LoRA, and MoReS), sampled from 100 instances each. Green represents visual representations, while grey indicates other (primarily textual) representations. Full fine-tuning and LoRA show strong reliance 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.

such as image understanding (Sidorov et al., 2020), visual question answering (Goyal et al., 2017a; Lu et al., 2022; Hudson & 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 fine-tune.

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 fine-tuning, 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 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; Shimamoto 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, multi-scale 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: To leverage pre-trained large language models (LLMs) for understanding visual instructions and generating responses, researchers have introduced cross-attention mechanisms to integrate image information into the language model. Notable examples include models such as LLaMA 3-V (Dubey et al., 2024), IDEFICS (Laureçon et al., 2023), and Flamingo (Awadalla et al., 2023; Alayrac et al., 2022). These models typically follow a two-stage training process: pretraining on large-scale image-text datasets, followed by supervised fine-tuning (SFT) with carefully curated high-quality data. During this process, the self-attention layers in the LLM decoder are kept frozen, with only the cross-attention and perceiver layers updated, ensuring that the text-only performance remains intact.

Another prominent approach employs a decoder-only architecture, as seen in models like the LLaVA family (Liu et al., 2024b;a; 2023), BLIP (Xue et al., 2024; Li et al., 2023a), and Qwen-VL (team, 2024; Bai et al., 2023). These models also follow the pretraining and visual instruction tuning paradigm. In the pretraining stage, a randomly initialized connector is trained while keeping the LLM frozen. However, recent studies (Bai et al., 2023; Chen et al., 2023) have demonstrated scenarios where both the projector and vision encoder are jointly trained during pretraining. Given the limited capacity of adapter modules, it is common to unfreeze the LLM during visual instruction tuning, while keeping the vision encoder frozen.

NVLM (Dai et al., 2024) represents a hybrid approach, combining elements of both the cross-attention and decoder-only architectures. In contrast, vision-encoder-free methods, as explored by models like Fuyu (Bavishi et al., 2023), SOLO (Chen et al., 2024b), and EVE (Diao et al., 2024), directly integrate visual information into LLMs at the pixel level, foregoing traditional vision encoders altogether.

While these approaches have advanced the integration of visual representations into LLMs, they still face significant challenges in the computational demands of visual instruction tuning, motivating further exploration into more efficient methods.

Visual Instruction Tuning: Fine tuning of multimodal large language models (MLLMs) for downstream tasks has gained considerable attention, but remains computationally expensive due to large-scale 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 & 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 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., 2022; 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} is the textual representation, R_{image} is the visual input representation, R_{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

216 fine-tuning, LoRA, and our proposed MoReS method. The results indicate that textual representa-
 217 tions often dominate the output generation process in both full fine-tuning and LoRA.
 218

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

222 To validate this observation, we introduce the Layer-wise Modality Attention Ratio (LMAR), for-
 223 mulated as:

$$224 \text{LMAR}_l = \frac{1}{N} \sum_{i=1}^N \frac{\alpha_l^{\text{image},i}}{\alpha_l^{\text{text},i}}, \quad (2)$$

225
 226
 227 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
 228 mean attention scores allocated to visual and textual tokens, respectively, in layer l for the i -th
 229 sample. LMAR thus provides a robust measure of the attention distribution between modalities,
 230 averaged over multiple samples to capture general trends in modality representation across layers.
 231

232
 233 In our experiments comparing various existing
 234 PEFT methods and full fine-tuning, IA3 (Liu
 235 et al., 2022) consistently achieves the highest
 236 average LMAR score across all layers while re-
 237 quiring the fewest trainable parameters. IA3’s
 238 superior performance can be attributed to its
 239 unique design, which introduces task-specific
 240 rescaling vectors that directly modulate key
 241 components of the Transformer architecture,
 242 such as the keys, values, and feed-forward lay-
 243 ers.

243 Unlike methods that introduce complex
 244 adapters or fine-tune all parameters, IA3 op-
 245 timizes a small but crucial set of parameters
 246 responsible for attention and representation
 247 learning. By applying element-wise scaling
 248 to the attention mechanisms, IA3 effectively
 249 re-balances the attention distribution across two
 250 intrinsic modalities. This design is particularly
 251 beneficial during visual instruction tuning, as
 252 it allows the model to dynamically reallocate
 253 more attention to visual representations without
 254 requiring many trainable parameters.

255 The identified relationship inspires that if the intrinsic modality imbalance can be addressed, the
 256 required number of trainable parameters can be potentially reduced further during visual instruction
 257 tuning. This offers a new direction for future improvements in PEFT methods for MLLMs.
 258

259 4 MORES METHOD

260
 261 Based on insights gained from intrinsic modality imbalance, we introduce Modality Linear
 262 Representation-Steering (**MoReS**) as a novel method for visual instruction tuning which can rebal-
 263 ance visual and textual representations and achieve comparable performance with fewer trainable
 264 parameters.
 265

266 Our approach is grounded in the linear subspace hypothesis, originally proposed by Bolukbasi et al.
 267 (2016), which suggests that information pertaining to a specific concept is encoded within a linear
 268 subspace in a model’s representation space. This hypothesis has been rigorously validated across
 269 numerous domains, including language understanding and interpretability (Lasri et al., 2022; Nanda
 et al., 2023; Amini et al., 2023; Wu et al., 2024c).

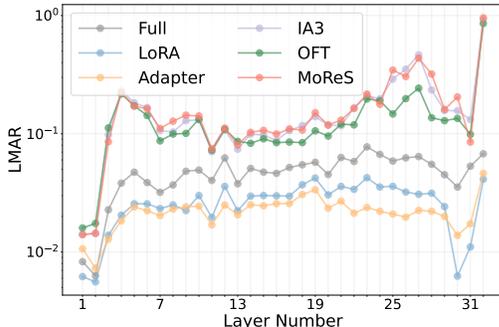


Figure 2: Layer-wise Modality Attention Ratio (LMAR) comparison across training methods, including Full fine-tuning, LoRA, Adapter, IA3, and our MoReS. Our MoReS method (red line) consistently demonstrates the highest LMAR across most layers, with a notable spike in the final layers. Compared with full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

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.

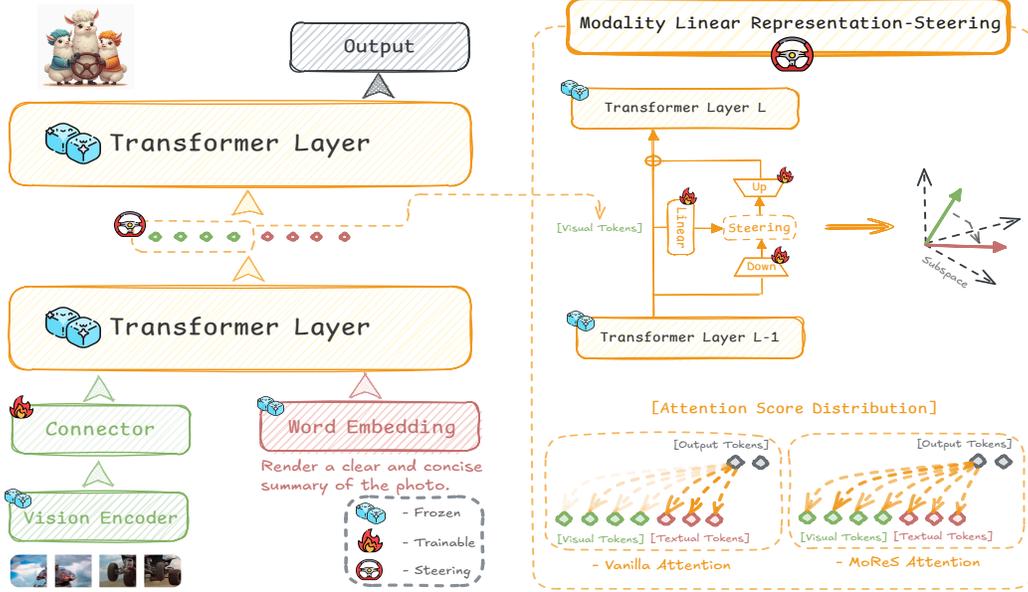


Figure 3: Schematic Overview of Modality Linear Representation-Steering (MoReS): **Left:** The architectural diagram depicts the integration of textual and visual tokens through transformer layers, leading to output token generation. **Right:** The mathematical formulation of MoReS illustrates the steering 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.

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:

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

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

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

where $W_{\text{down}} \in \mathbb{R}^{d \times D}$ is a downsampling matrix. To preserve the fidelity of the representation and 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 5.5.

In Section A.1, 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 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.

5.1 EXPERIMENT SETTINGS

5.1.1 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.

5.1.2 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 $\mathbf{W}_{\text{up}} \in \mathbb{R}^{d \times r}$, where d is the hidden layer dimension and r is the bottleneck dimension. The adapter output is computed as:

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

where the residual connection ($+\mathbf{x}$) 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:

$$W_0 + \Delta W = W_0 + BA \quad (6)$$

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:

$$z = W^\top x = (R \cdot W_0)^\top x \quad (7)$$

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:

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

where \odot denotes multiplication by element.

5.2 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 multiple benchmarks, including VQAv2 (Goyal et al., 2017b) and GQA (Hudson & Manning, 2019), which test visual perception through open-ended short answers, and VizWiz (Gurari et al., 2018), with 8,000 images designed for zero-shot generalization in visual questions posed by visually impaired individuals. We also use the image subset of ScienceQA (Lu et al., 2022) with multiple-choice questions to assess zero-shot scientific question answering, while TextVQA (Singh et al., 2019) measures performance on text-rich visual questions. MM-Vet (Yu et al., 2023) evaluates the model’s ability to engage in visual conversations, with correctness and helpfulness scored by GPT-4. Additionally, POPE (Li et al., 2023b) quantifies hallucination of MLLMs. Finally, we apply the MMMU benchmark (Yue et al., 2024) to assess core multimodal skills, including perception, knowledge, and reasoning.

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.

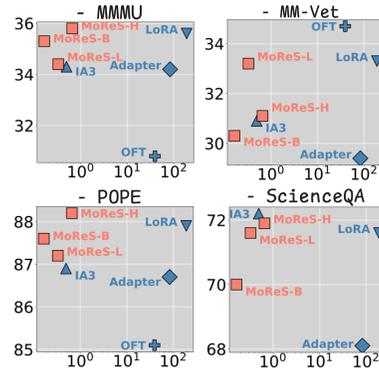


Figure 4: Comparison of parameter count vs. performance for MoReS and other PEFT methods across four benchmarks.

Model	Method	TP*	VQAv2	GQA	TextVQA	SciQA-IMG	POPE	MM-Vet	MMMU	Avg
LLaVA Steering-3B	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
	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

5.3 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.

Model	Method	TP*	SciQA-IMG	VizWiz	IconQA-txt	IconQA-blank	Scale	Method	TP*	SciQA-IMG	VizWiz	IconQA
LLaVA Steering-3B	Adapter	83M	92.3	62.9	93.5	95.8	Small	FT	2.78B	33.8	51.2	68.1
	LoRA	188.7M	93.9	61.6	93.9	96.5		Adapter	83M	81.0	57.4	72.4
	OFT	39.32M	86.3	42.0	87.8	42.0		LoRA	188.74M	84.0	58.5	74.2
	IA3	0.492M	90.2	58.4	84.5	94.7		OFT	39.32M	79.2	43.2	35.9
	MoReS-B	0.164M	89.7	59.2	84.0	94.2		IA3	0.492M	79.9	50.5	73.0
LLaVA Steering-7B	Adapter	201.3M	82.7	59.7	72.1	71.6	Medium	MoReS-L	0.328M	78.2	55.0	69.7
	LoRA	319.8M	87.6	60.6	77.7	70.2		FT	2.78B	78.2	58.9	92.2
	OFT	100.7M	78.3	55.1	19.4	22.7		Adapter	83M	92.1	60.6	93.2
	IA3	0.614M	83.8	54.3	65.1	70.4		LoRA	188.74M	92.9	60.5	92.7
	MoReS-B	0.262M	83.6	54.2	64.2	70.2		OFT	39.32M	86.4	44.4	45.5
LLaVA Steering-13B	Adapter	314.6M	87.9	61.4	78.2	73.0	Large	IA3	0.492M	91.9	57.1	90.6
	LoRA	500.7M	92.1	62.0	80.2	73.2		MoReS-L	0.328M	92.1	56.6	89.9
	OFT	196.6M	82.7	59.5	3.4	22.3		FT	2.78B	88.9	59.4	95.7
	IA3	0.963M	90.5	54.6	73.8	71.7		Adapter	83M	92.4	61.3	95.2
	MoReS-B	0.410M	89.5	54.3	74.9	71.5		LoRA	188.74M	93.9	61.8	96.0

Table 2: Results of Task-Specific Fine-tuning, where higher values correspond to better performance.

Table 3: Results of multi-scale tasks.

5.4 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, 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 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.

5.5 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 4 summarizes 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 overhead due to its sparse steering approach.

6 LLAVA STEERING FACTORY

We identified a pressing need for a comprehensive framework to systematically analyze and compare various model architectures and training strategies in MLLMs. The diversity of design choices and

Subspace Rank	TP*	SciQA-IMG	VizWiz	IconQA-txt	IconQA-blank	Avg	Steered Visual Token Ratio	SciQA-IMG	VizWiz	IconQA-txt	IconQA-blank
1	0.164M	89.6	59.2	84.0	94.2	81.8	1%	89.7	59.2	84.0	94.1
2	0.328M	89.7	59.2	83.9	94.0	81.7	25%	89.9	59.0	80.2	93.8
4	0.655M	89.5	58.7	83.8	94.1	81.5	50%	88.9	59.0	79.8	92.6
8	1.340M	89.6	58.9	83.7	93.9	81.5	100%	85.8	60.5	67.7	87.8

Table 4: Results of the subspace rank choice and steered visual token ratio. The grey shading indicates the best results and our selected parameters.

techniques complicates the standardization and understanding of how these components interact. Evaluating each method across different open-source models is often time-consuming and inconsistent due to implementation differences, which necessitate extensive data preprocessing and careful alignment between architectures and training recipes.

In the LLaVA Steering Factory, we establish standardized training and evaluation pipelines, along with flexible data preprocessing and model configurations. Our framework allows researchers to easily customize their models with various training strategies without the need for additional coding. We implement all mainstream LLMs and vision encoders, including multiple PEFT methods and our proposed MoReS technique. Furthermore, we support a wide range of benchmarks and integrate our intrinsic modality imbalance evaluation. The goal of the LLaVA Steering Factory is to facilitate research in MLLMs, particularly in addressing intrinsic modality imbalance to optimize visual instruction tuning.

An overview of the main components of the LLaVA Steering Factory is provided in Figure 5.

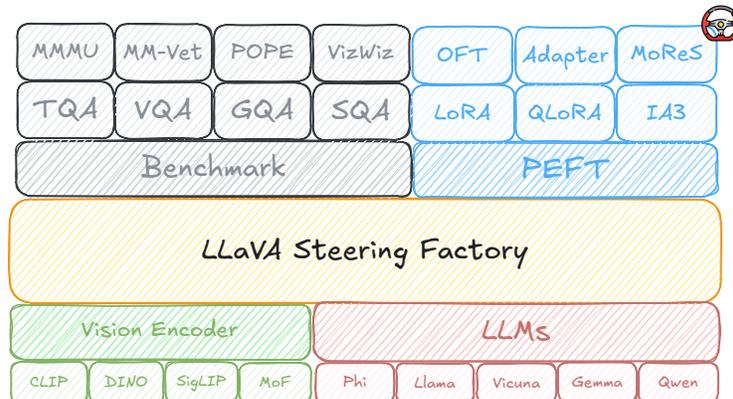


Figure 5: Architectural overview of the proposed LLaVA Steering Factory: A Modular Codebase for MLLMs.

7 CONCLUSION

This paper introduces Modality Linear Representation-Steering (**MoReS**), a novel method to significantly reduce the required number of trainable parameters during visual instruction tuning. The key idea behind MoReS is to re-balance visual and textual representations while still maintaining strong performance across a variety of downstream tasks. By integrating MoReS into LLaVA family models, comprehensive evaluation results confirm the effectiveness of the proposed solution. Hence, it further confirms our assertion that intrinsic modality rebalance would represent a promising new approach to optimizing visual instruction tuning.

To facilitate future research in the community, we also present the LLaVA Steering Factory, a versatile framework designed to enhance the development and evaluation of MLLMs with minimal coding effort. This framework enables customizable training configurations for various tasks and integrates analytical tools, such as attention score distribution analysis. This facilitates systematic comparisons among different methods and offers deeper insights into the intrinsic modality imbalance, ultimately contributing to more effective visual instruction tuning.

REFERENCES

- 540
541
542 Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen
543 Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko,
544 Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dong-
545 dong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang
546 Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit
547 Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao,
548 Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin
549 Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim,
550 Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden,
551 Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong
552 Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro
553 Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norrick, Barun Patra, Daniel Perez-
554 Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo
555 de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim,
556 Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla,
557 Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua
558 Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp
559 Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Ji-
560 long Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan,
561 Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan
562 Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your
563 phone, 2024. URL <https://arxiv.org/abs/2404.14219>.
- 564
565 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson,
566 Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza
567 Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Mon-
568 teiro, Jacob L Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Shar-
569 ifzadeh, Miłkoł aj Bińkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén
570 Simonyan. Flamingo: a visual language model for few-shot learning. In S. Koyejo,
571 S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neu-
572 ral Information Processing Systems*, volume 35, pp. 23716–23736. Curran Associates, Inc.,
573 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/
574 file/960a172bc7fbf0177cccbb411a7d800-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/960a172bc7fbf0177cccbb411a7d800-Paper-Conference.pdf).
- 575
576 Afra Amini, Tiago Pimentel, Clara Meister, and Ryan Cotterell. Naturalistic causal probing for
577 morpho-syntax. *Transactions of the Association for Computational Linguistics*, 11:384–403,
578 2023.
- 579
580 Matan Avitan, Ryan Cotterell, Yoav Goldberg, and Shauli Ravfogel. Natural language counterfactu-
581 als through representation surgery, 2024. URL <https://arxiv.org/abs/2402.11355>.
- 582
583 Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani
584 Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei
585 Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo: An open-
586 source framework for training large autoregressive vision-language models. *arXiv preprint
587 arXiv:2308.01390*, 2023.
- 588
589 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang
590 Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, local-
591 ization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- 592
593 Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani,
594 and Sağnak Taşırılar. Introducing our multimodal models, 2023. URL [https://www.adept.
595 ai/blog/fuyu-8b](https://www.adept.ai/blog/fuyu-8b).
- 596
597 Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is
598 to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in
599 neural information processing systems*, 29, 2016.

- 594 Guanzheng Chen, Fangyu Liu, Zaiqiao Meng, and Shangsong Liang. Revisiting parameter-efficient
595 tuning: Are we really there yet? In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.),
596 *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp.
597 2612–2626, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational
598 Linguistics. doi: 10.18653/v1/2022.emnlp-main.168. URL [https://aclanthology.org/
599 2022.emnlp-main.168](https://aclanthology.org/2022.emnlp-main.168).
- 600 Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang.
601 An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-
602 language models. *arXiv preprint arXiv:2403.06764*, 2024a.
- 603 Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua
604 Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint
605 arXiv:2311.12793*, 2023.
- 606 Yangyi Chen, Xingyao Wang, Hao Peng, and Heng Ji. A single transformer for scalable vision-
607 language modeling. *arXiv preprint arXiv:2407.06438*, 2024b.
- 608 Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong,
609 Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to com-
610 mercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024c.
- 611 Wenliang Dai, Nayeon Lee, Boxin Wang, Zhuoling Yang, Zihan Liu, Jon Barker, Tuomas Rinta-
612 maki, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. Nvlm: Open frontier-class multi-
613 modal llms, 2024. URL <https://arxiv.org/abs/2409.11402>.
- 614 Haiwen Diao, Yufeng Cui, Xiaotong Li, Yueze Wang, Huchuan Lu, and Xinlong Wang. Unveiling
615 encoder-free vision-language models. *arXiv preprint arXiv:2406.11832*, 2024.
- 616 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
617 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony
618 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,
619 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,
620 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris
621 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,
622 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny
623 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,
624 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael
625 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-
626 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah
627 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan
628 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-
629 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy
630 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,
631 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-
632 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,
633 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearly, Laurens van der
634 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,
635 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-
636 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,
637 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,
638 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur
639 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-
640 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,
641 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,
642 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-
643 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,
644 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,
645 Sharath Rapparthi, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,
646 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney
647 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,

648 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,
649 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-
650 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,
651 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur,
652 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre
653 Couderc, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha
654 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay
655 Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda
656 Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew
657 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita
658 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh
659 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De
660 Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-
661 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina
662 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,
663 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,
664 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana
665 Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,
666 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-
667 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco
668 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella
669 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory
670 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,
671 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-
672 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman,
673 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer
674 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe
675 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie
676 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun
677 Zand, Kathy Matosich, Kaushik Veeraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal
678 Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,
679 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian
680 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,
681 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-
682 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel
683 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-
684 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-
685 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,
686 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,
687 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,
688 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,
689 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,
690 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,
691 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-
692 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-
693 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang
694 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen
695 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho,
696 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser,
697 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-
698 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,
699 Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu
700 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-
701 stable, Xiaocheng Tang, Xiaofang Wang, Xiaoqian Wu, Xiaolan Wang, Xide Xia, Xilun Wu,
Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,
Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef
Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024.
URL <https://arxiv.org/abs/2407.21783>.

- 702 Atticus Geiger, Zhengxuan Wu, Christopher Potts, Thomas Icard, and Noah Goodman. Find-
703 ing alignments between interpretable causal variables and distributed neural representations. In
704 *Causal Learning and Reasoning*, pp. 160–187. PMLR, 2024.
- 705 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V
706 in VQA matter: Elevating the role of image understanding in Visual Question Answering. In
707 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017a.
- 708 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa
709 matter: Elevating the role of image understanding in visual question answering. In *Proceedings*
710 *of the IEEE conference on computer vision and pattern recognition*, pp. 6904–6913, 2017b.
- 711 Clément Guerner, Anej Svete, Tianyu Liu, Alexander Warstadt, and Ryan Cotterell. A geometric
712 notion of causal probing. *arXiv preprint arXiv:2307.15054*, 2023.
- 713 Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and
714 Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In
715 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3608–3617,
716 2018.
- 717 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, An-
718 drea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp.
719 In *International conference on machine learning*, pp. 2790–2799. PMLR, 2019.
- 720 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,
721 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
722 *arXiv:2106.09685*, 2021.
- 723 Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning
724 and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer*
725 *vision and pattern recognition*, pp. 6700–6709, 2019.
- 726 Karim Lasri, Tiago Pimentel, Alessandro Lenci, Thierry Poibeau, and Ryan Cotterell. Probing for
727 the usage of grammatical number. *arXiv preprint arXiv:2204.08831*, 2022.
- 728 Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov,
729 Thomas Wang, Siddharth Karamcheti, Alexander M. Rush, Douwe Kiela, Matthieu Cord, and
730 Victor Sanh. Obelics: An open web-scale filtered dataset of interleaved image-text documents,
731 2023. URL <https://arxiv.org/abs/2306.16527>.
- 732 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: bootstrapping language-image
733 pre-training with frozen image encoders and large language models. In *Proceedings of the 40th*
734 *International Conference on Machine Learning, ICML’23*. JMLR.org, 2023a.
- 735 Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time
736 intervention: Eliciting truthful answers from a language model. *Advances in Neural Information*
737 *Processing Systems*, 36, 2024.
- 738 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.
739 In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th*
740 *Annual Meeting of the Association for Computational Linguistics and the 11th International Joint*
741 *Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 4582–4597, Online,
742 August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.353.
743 URL <https://aclanthology.org/2021.acl-long.353>.
- 744 Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating
745 object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023b.
- 746 Yanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee.
747 Textbooks are all you need ii: **phi-1.5** technical report. *arXiv preprint arXiv:2309.05463*, 2023c.
- 748 Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and
749 Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context
750 learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022.

- 756 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In
757 A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in*
758 *Neural Information Processing Systems*, volume 36, pp. 34892–34916. Curran Associates, Inc.,
759 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/](https://proceedings.neurips.cc/paper_files/paper/2023/file/6dcf277ea32ce3288914faf369fe6de0-Paper-Conference.pdf)
760 [file/6dcf277ea32ce3288914faf369fe6de0-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/6dcf277ea32ce3288914faf369fe6de0-Paper-Conference.pdf).
- 761 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction
762 tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
763 *(CVPR)*, pp. 26296–26306, June 2024a.
- 764 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.
765 Llava-next: Improved reasoning, ocr, and world knowledge, January 2024b. URL [https://](https://llava-vl.github.io/blog/2024-01-30-llava-next/)
766 llava-vl.github.io/blog/2024-01-30-llava-next/.
- 767 Zikang Liu, Kun Zhou, Wayne Xin Zhao, Dawei Gao, Yaliang Li, and Ji-Rong Wen. Less is more:
768 Data value estimation for visual instruction tuning. *arXiv preprint arXiv:2403.09559*, 2024c.
- 769 Pan Lu, Liang Qiu, Jiaqi Chen, Tony Xia, Yizhou Zhao, Wei Zhang, Zhou Yu, Xiaodan Liang,
770 and Song-Chun Zhu. Iconqa: A new benchmark for abstract diagram understanding and visual
771 language reasoning. *arXiv preprint arXiv:2110.13214*, 2021.
- 772 Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,
773 Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for
774 science question answering. In *The 36th Conference on Neural Information Processing Systems*
775 *(NeurIPS)*, 2022.
- 776 Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter,
777 Dhruvi Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. Mm1: Methods, analysis & insights
778 from multimodal llm pre-training. *arXiv preprint arXiv:2403.09611*, 2024.
- 779 Neel Nanda, Andrew Lee, and Martin Wattenberg. Emergent linear representations in world models
780 of self-supervised sequence models. *arXiv preprint arXiv:2309.00941*, 2023.
- 781 Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian Weller,
782 and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetuning. *Advances*
783 *in Neural Information Processing Systems*, 36:79320–79362, 2023.
- 784 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
785 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
786 models from natural language supervision. In *International conference on machine learning*, pp.
787 8748–8763. PMLR, 2021.
- 788 Erica K. Shimomoto, Edison Marrese-Taylor, Hiroya Takamura, Ichiro Kobayashi, and Yusuke
789 Miyao. A subspace-based analysis of structured and unstructured representations in image-text
790 retrieval. In Wenjuan Han, Zilong Zheng, Zhouhan Lin, Lifeng Jin, Yikang Shen, Yoon Kim,
791 and Kewei Tu (eds.), *Proceedings of the Workshop on Unimodal and Multimodal Induction of*
792 *Linguistic Structures (UM-IoS)*, pp. 29–44, Abu Dhabi, United Arab Emirates (Hybrid), Decem-
793 ber 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.umios-1.4. URL
794 <https://aclanthology.org/2022.umios-1.4>.
- 795 Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for
796 image captioning with reading comprehension. 2020.
- 797 Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh,
798 and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF*
799 *conference on computer vision and pattern recognition*, pp. 8317–8326, 2019.
- 800 Shashwat Singh, Shauli Ravfogel, Jonathan Herzig, Roei Aharoni, Ryan Cotterell, and Ponnur-
801 rangam Kumaraguru. Representation surgery: Theory and practice of affine steering, 2024. URL
802 <https://arxiv.org/abs/2402.09631>.
- 803 Nishant Subramani, Nivedita Suresh, and Matthew E Peters. Extracting latent steering vectors from
804 pretrained language models. *arXiv preprint arXiv:2205.05124*, 2022.

- 810 Qwen team. Qwen2-vl. 2024.
811
- 812 Yaqing Wang, Sahaj Agarwal, Subhabrata Mukherjee, Xiaodong Liu, Jing Gao, Ahmed Has-
813 san Awadallah, and Jianfeng Gao. AdaMix: Mixture-of-adaptations for parameter-efficient
814 model tuning. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of*
815 *the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 5744–5760,
816 Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguis-
817 tics. doi: 10.18653/v1/2022.emnlp-main.388. URL [https://aclanthology.org/2022.](https://aclanthology.org/2022.emnlp-main.388)
818 [emnlp-main.388](https://aclanthology.org/2022.emnlp-main.388).
- 819 Muling Wu, Wenhao Liu, Xiaohua Wang, Tianlong Li, Changze Lv, Zixuan Ling, Jianhao Zhu,
820 Cenyuan Zhang, Xiaoqing Zheng, and Xuanjing Huang. Advancing parameter efficiency in fine-
821 tuning via representation editing. *arXiv preprint arXiv:2402.15179*, 2024a.
- 822
823 Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D. Man-
824 ning, and Christopher Potts. Reft: Representation finetuning for language models, 2024b. URL
825 <https://arxiv.org/abs/2404.03592>.
- 826
827 Zhengxuan Wu, Atticus Geiger, Thomas Icard, Christopher Potts, and Noah Goodman. Interpretabil-
828 ity at scale: Identifying causal mechanisms in alpaca. *Advances in Neural Information Processing*
829 *Systems*, 36, 2024c.
- 830
831 Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Vi-
832 raj Prabhu, Yutong Dai, Michael S Ryoo, Shrikant Kendre, Jieyu Zhang, Can Qin, Shu Zhang,
833 Chia-Chih Chen, Ning Yu, Juntao Tan, Tulika Manoj Awalganekar, Shelby Heinecke, Huan
834 Wang, Yejin Choi, Ludwig Schmidt, Zeyuan Chen, Silvio Savarese, Juan Carlos Niebles, Caim-
835 ing Xiong, and Ran Xu. xgen-mm (blip-3): A family of open large multimodal models, 2024.
URL <https://arxiv.org/abs/2408.08872>.
- 836
837 Linli Yao, Lei Li, Shuhuai Ren, Lean Wang, Yuanxin Liu, Xu Sun, and Lu Hou. Deco: Decou-
838 pling token compression from semantic abstraction in multimodal large language models. *arXiv*
839 *preprint arXiv:2405.20985*, 2024.
- 840
841 Ting Yao, Yingwei Pan, Chong-Wah Ngo, Houqiang Li, and Tao Mei. Semi-supervised domain
842 adaptation with subspace learning for visual recognition. In *Proceedings of the IEEE Conference*
843 *on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- 844
845 Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang,
846 and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv*
847 *preprint arXiv:2308.02490*, 2023.
- 848
849 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens,
850 Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multi-
851 modal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF*
852 *Conference on Computer Vision and Pattern Recognition*, pp. 9556–9567, 2024.
- 853
854 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language
855 image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer*
856 *Vision*, pp. 11975–11986, 2023.
- 857
858 Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small
859 language model, 2024. URL <https://arxiv.org/abs/2401.02385>.
- 860
861 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
862 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Sto-
863 ica. Judging llm-as-a-judge with mt-bench and chatbot arena. In A. Oh, T. Naumann,
A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Informa-*
tion Processing Systems, volume 36, pp. 46595–46623. Curran Associates, Inc., 2023a.
URL [https://proceedings.neurips.cc/paper_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/91f18a1287b398d378ef22505bf41832-Paper-Datasets_and_Benchmarks.pdf)
[91f18a1287b398d378ef22505bf41832-Paper-Datasets_and_Benchmarks.](https://proceedings.neurips.cc/paper_files/paper/2023/file/91f18a1287b398d378ef22505bf41832-Paper-Datasets_and_Benchmarks.pdf)
pdf.

864 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
865 Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
866 Judging llm-as-a-judge with mt-bench and chatbot arena, 2023b.
867

868 Baichuan Zhou, Ying Hu, Xi Weng, Junlong Jia, Jie Luo, Xien Liu, Ji Wu, and Lei Huang. Tynllava:
869 A framework of small-scale large multimodal models, 2024a. URL <https://arxiv.org/abs/2402.14289>.
870

871 Xiongtao Zhou, Jie He, Yuhua Ke, Guangyao Zhu, Víctor Gutiérrez-Basulto, and Jeff Z Pan. An
872 empirical study on parameter-efficient fine-tuning for multimodal large language models. *arXiv*
873 *preprint arXiv:2406.05130*, 2024b.
874

875 Xingyu Zhu, Beier Zhu, Yi Tan, Shuo Wang, Yanbin Hao, and Hanwang Zhang. Selective vision-
876 language subspace projection for few-shot clip. *arXiv preprint arXiv:2407.16977*, 2024.

877 Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan,
878 Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engineering: A
879 top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*, 2023.
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

A APPENDIX

A.1 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:

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

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

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

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

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

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:

$$\text{MoReS}(h) = W_{\text{up}} \cdot \phi(h), \quad \text{where} \quad \phi(h) = \text{Linear}(h) - W_{\text{down}} h \quad (12)$$

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

$$A'_{i,\text{image}} = \text{softmax} \left(\frac{(h_i W_q)(h'_{\text{image}} W_k)^T}{\sqrt{D}} \right) \quad (13)$$

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:

$$\|\Delta y\|_2 = \|W_o(C'_{\text{image}} - C_{\text{image}})\|_2 \leq \|W_o\|_2 \|C'_{\text{image}} - C_{\text{image}}\|_2 \quad (14)$$

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_{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'_{image} .

To ensure stability while allowing for this significant attention redistribution, we consider the Lipschitz continuity of the model:

$$\|f(h'_{\text{image}}) - f(h_{\text{image}})\|_2 \leq L \|h'_{\text{image}} - h_{\text{image}}\|_2 \quad (15)$$

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.2 IMPLEMENTATION DETAIL

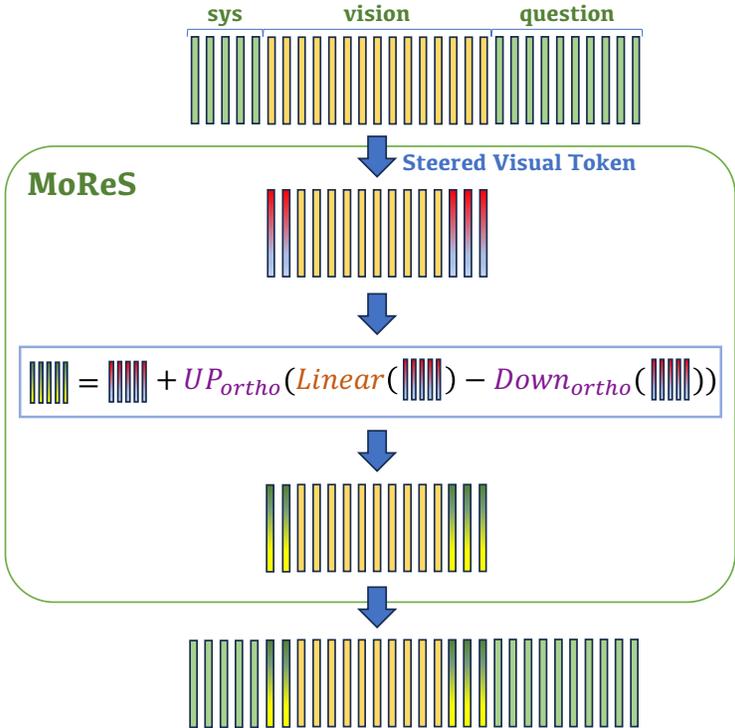


Figure 6: MoReS module flowchart.

Regarding the implementation, we have adopted a highly modular design for the LLM, integrating it with MoReS to enable precise steering at specific 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 6) and an UML diagram (Figure 7) here, which detail the implementation process.

A.3 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

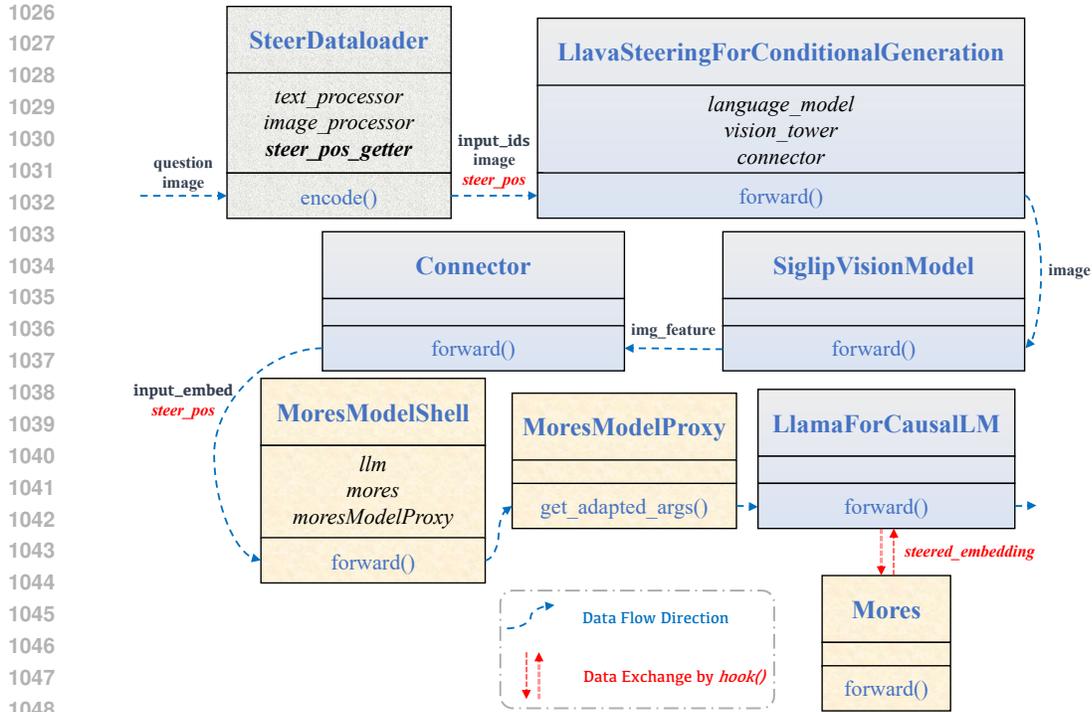


Figure 7: The UML diagram for MoReS

avoided, leading to faster and more efficient inference while maintaining the essential information needed for accurate predictions.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

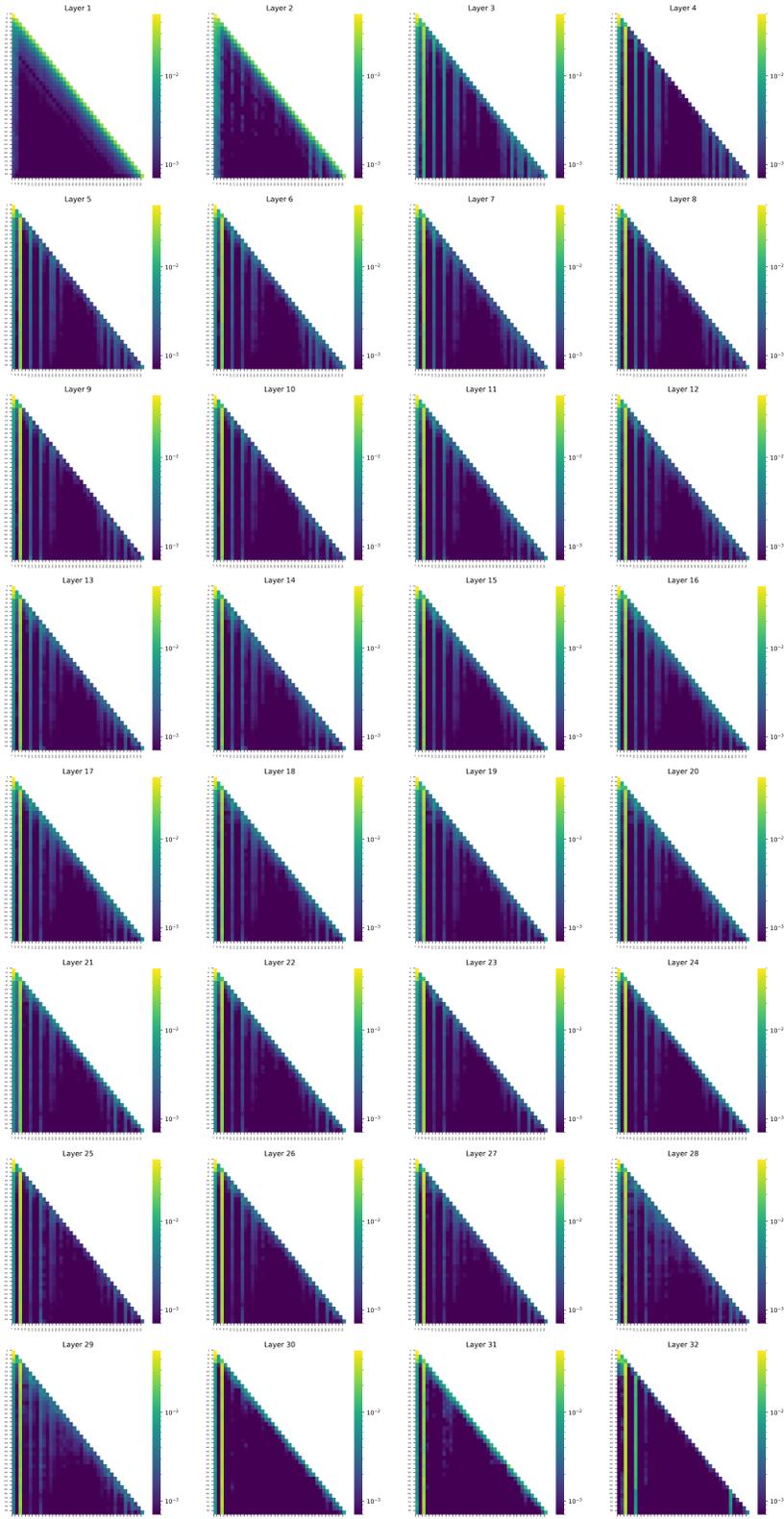


Figure 8: Full Attention Maps of Each Layer