

Visual Representation Alignment for Multimodal Large Language Models

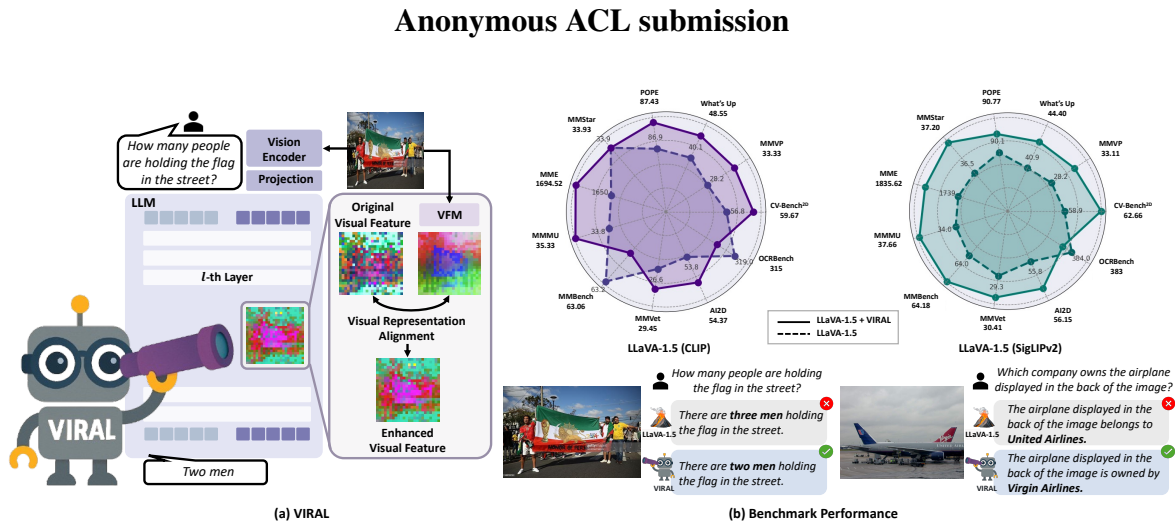


Figure 1: (a) **Visual Representation Alignment (VIRAL)** introduces an auxiliary regularization objective on the visual pathway, preventing MLLMs from discarding detailed attributes of the input vision encoder during training while incorporating additional visual knowledge from vision foundation models (VFMs). (b) When trained with DINOv2 (Oquab et al., 2023) as the VFM, VIRAL consistently yields more accurate visually grounded responses and achieves substantial improvements over standard baselines (Liu et al., 2023) across diverse vision encoders, including CLIP (Radford et al., 2021) and SigLIPv2 (Tschannen et al., 2025).

Abstract

Multimodal large language models (MLLMs) trained with visual instruction tuning have achieved strong performance across diverse tasks, yet they remain limited in vision-centric tasks such as object counting and spatial reasoning. We find that this limitation arises not merely from the choice of vision encoder, but from the lack of explicit supervision on visual representations during training, which causes detailed visual information to be gradually weakened even when strong vision foundation models (VFMs) are used as vision encoders. To this end, we present **Visual Representation Alignment (VIRAL)**, a simple yet effective regularization strategy that aligns the internal visual representations of MLLMs with those of pre-trained VFMs. By explicitly enforcing this alignment, VIRAL preserves rich visual information within the MLLM while enabling it to leverage complementary visual knowledge from VFMs, thereby enhancing its ability to reason over complex visual inputs.

1 Introduction

Recent advancements in multimodal large language models (MLLMs) (OpenAI, 2023; Bai et al., 2023a; Team et al., 2023; Chen et al., 2024c), particularly those employing visual instruction tuning (Liu

et al., 2023), have achieved notable success in diverse multimodal tasks. By connecting pretrained large language models (LLMs) (Touvron et al., 2023; Chiang et al., 2023; Chen et al., 2024c; Bai et al., 2025) with vision encoders (Radford et al., 2021; Chen et al., 2024c; Tong et al., 2024a) through a lightweight vision–language projector, visual instruction tuning enables LLMs to effectively interpret and reason over visual inputs.

Despite these successes, numerous studies report persistent limitations in vision-centric tasks such as object counting and spatial reasoning (Tong et al., 2024b; Qi et al., 2025; Yuksekgonul et al., 2022; Ma et al., 2023). Early approaches largely attribute these shortcomings to the visual encoder or the vision-language projector. In response, subsequent works have introduced stronger vision encoders (Lu et al., 2024; Li et al., 2024) and more expressive projectors (Liu et al., 2024a; Cha et al., 2024; McKinzie et al., 2024), aiming to supply the language model with richer and more comprehensive visual representations. While they yield notable improvements, incorporating powerful vision encoders or projectors are inherently constrained in scalability, with recent works (Fu et al., 2025) suggesting that limitations in vision-centric tasks come from failure of the LLMs to effectively utilize available

visual information, rather than the quality of the visual encoder’s representation itself.

These prior findings motivate us to seek further advances beyond architectural refinement alone. In this paper, we first revisit the conventional training paradigm of visual instruction tuning. Existing MLLMs are predominantly fine-tuned with a language-modeling objective, updating both the LLM and the vision-language projector while concentrating supervision almost entirely on textual outputs (Li et al., 2024; Bai et al., 2023b; Chen et al., 2024c). As a result, visual tokens receive only indirect, language-mediated supervision despite comprising a substantial fraction of the multimodal input. In effect, the visual pathway remains under-supervised, raising a central question: *Apart from the vision encoder’s representation quality itself, is the prevailing multimodal training setup adequate for capturing and preserving visual information?*

We hypothesize that text-only supervision encourages the model to retain only those visual details that immediately aid text prediction, discarding other potentially useful cues. For instance, as in examples shown in Fig. 1, a caption such as “A photo of a group of people holding a large flag.” provides little incentive to preserve the flag’s color, the exact number of people, or their spatial layout—attributes needed for fine-grained image understanding. In short, text-only supervision aligns visual features with language efficiently, but does so at the cost of losing the richer and more structured representations provided by the vision encoder (Venhoff et al., 2025; Neo et al., 2024).

To validate this hypothesis, we conduct an experiment (see Fig. 2) and observe that visual representations within MLLMs trained under exclusive textual supervision rapidly diverge from those produced by the input vision encoder, which we refer to as *visual representation misalignment*. In addition, we observe that simply enforcing alignment of the internal visual representations of MLLMs (features that lost detailed visual information cues) with the model’s own vision encoder (features with richer visual cues) to compensate for the lost visual information, already yields substantial gains in fine-grained visual understanding, indicating that MLLMs discard visual information during training and retaining them is important for fine-grained image understandings.

Motivated by these findings, we propose **Visual Representation Alignment (VIRAL)**, a simple

yet effective regularization strategy that directly supervises the visual pathway in MLLMs by aligning their internal visual representations to a reference vision model. While utilizing the model’s own vision encoder as alignment reference effectively improves performance in vision-centric tasks, we further find that substantially larger gains are achieved when stronger vision foundation models (VFMs) (Oquab et al., 2023; Kirillov et al., 2023; Yang et al., 2024; Ranzinger et al., 2024) are used instead. Trained with vision-centric objectives, VFMs act as a stronger alignment reference that not only prevents visual information loss but also provides complementary visual knowledge absent from the model’s own vision encoder. Through extensive experiments on various multimodal benchmarks, we demonstrate that VIRAL consistently delivers significant improvements across all tasks.

We summarize our contributions as follows:

- We show that, under the visual instruction tuning paradigm, internal visual representations in MLLMs often drift from those of their vision encoders, leading to degraded spatial reasoning due to the loss of fine-grained visual information.
- We propose **VIRAL**, a novel regularization strategy that explicitly aligns MLLM visual representations with pretrained VFM features, preventing the loss of fine-grained attributes and enabling richer visual understanding.
- We show consistent improvements of an average **9.4%** over the baseline on vision-centric benchmarks, with extensive ablation studies and analysis to validate our design choices.

2 Related Work

Visual information processing in MLLMs. Recent studies (Kaduri et al., 2025; Zhang et al., 2025b; Jiang et al., 2025; Kang et al., 2025a) have begun to analyze how visual information is processed and transformed inside MLLMs. Prior work reveals a structured hierarchy in MLLMs, where early layers encode global visual context, intermediate layers capture fine-grained spatial information, and later layers integrate multimodal signals for generation (Kaduri et al., 2025; Zhang et al., 2025b). Within this hierarchy, the intermediate layers have been shown to play a critical role in visual grounding and spatial reasoning.

Complementary to these structural analyses, recent work has highlighted limitations in how MLLMs utilize visual representations internally. Fu et al. (2025) show that even when high-quality visual representations are available within the model, the language model often fails to effectively exploit them for vision-centric tasks. Consistent with these observations, our analysis shows that fine-grained visual information is progressively lost during training, and that preserving such information in the intermediate layers—where spatial semantics emerge—is particularly important for vision-centric reasoning.

Improving visual information in MLLMs.

While recent works have increasingly examined the internal information flow of MLLMs, most prior efforts remain concentrated on the input stage—particularly the use of frozen vision encoders. Improvements at this stage have largely focused on adopting stronger or multiple vision encoders (Kar et al., 2024; Lu et al., 2024; Shi et al., 2024; Azadani et al., 2025) or enhancing efficiency by reducing the overhead of visual tokens (Vasu et al., 2025; Yang et al., 2025; Wen et al., 2025). These advances have proven valuable, yet they primarily address the quality and efficiency of the initial visual representations, with comparatively less attention given to how visual information is processed and lost once injected into the model.

Representation supervisions. Representation-level supervision has recently gained attention in generative modeling. Methods such as REPA (Yu et al., 2024) align internal representations with those of strong vision encoders to improve learnability during training (Labs, 2025). In MLLMs, several recent works (Wang et al., 2024; Jain et al., 2025) similarly explore representation supervision to enhance visual perception. However, these methods simply utilize representations from VFMs to aid MLLMs to better solve downstream tasks such as providing geometric foundation features for better geometry understandings (Jain et al., 2025). On the other hand, our work is motivated by analyzing why current MLLMs fail in fine-grained image understanding, where we identify a key problem of visual information loss and introduce representation alignment as an auxiliary task to prevent such visual information loss for improved capabilities in downstream vision-centric tasks.

3 Preliminaries

Multimodal large language models (MLLMs).

MLLMs typically consist of a pre-trained language model $LM_\theta(\cdot)$ and a vision encoder $V_\psi(\cdot)$, which is connected with a vision-language projector $P_\phi(\cdot)$, where θ , ψ , and ϕ denote corresponding learnable parameters. To generate answers grounded on both input image and text, the frozen vision encoder $V_\psi(\cdot)$ first extracts patch-level features from an input image $I \in \mathbb{R}^{H \times W \times 3}$ with height H and width W such that $\mathbf{z} = V_\psi(I) \in \mathbb{R}^{N \times D_z}$, where N and D_z denote the number of visual tokens and the dimension of the visual features, respectively. In typical MLLMs, a linear vision-language projector $P_\phi(\cdot)$ maps these features into the language model’s embedding space, producing a sequence of visual tokens $\mathbf{e}^{\text{img}} = P_\phi(\mathbf{z}) \in \mathbb{R}^{N \times D}$, where D is the hidden dimension of the language model. Text inputs are embedded into the \mathbf{h} that $\mathbf{e}^{\text{text}} \in \mathbb{R}^{K \times D}$, where K denotes the length of the text tokens. and the language model processes the concatenated multimodal sequence.

During training, MLLMs are typically optimized using a text-only language modeling objective, where supervision is applied solely to the output text tokens:

$$\mathcal{L}_{\text{LM}} = -\frac{1}{K} \sum_{i=1}^K \log p_{\theta, \phi}(\mathbf{e}_i^{\text{text}} | \mathbf{e}_{<i}^{\text{text}}, \mathbf{e}^{\text{img}}). \quad (1)$$

As a result, visual representations \mathbf{e}^{img} receive no explicit vision-specific supervision, and all learning signals are mediated through language prediction.

4 Methodology

4.1 Do MLLMs lose visual information?

While MLLMs take a substantial number of visual tokens as input, they are typically trained with a text-only language modeling loss applied to the output text tokens. Consequently, all learning signals are mediated through language supervision, and the visual representations \mathbf{e}^{img} receive no vision-specific supervision, as illustrated in Fig. 2-(a). In the absence of explicit visual supervision, we hypothesize that the model learns to prioritize only those visual features that immediately aid textual prediction, often discarding other potentially useful information. This, in turn, causes the internal visual representations to drift away from the rich features produced by the vision encoder—an effect

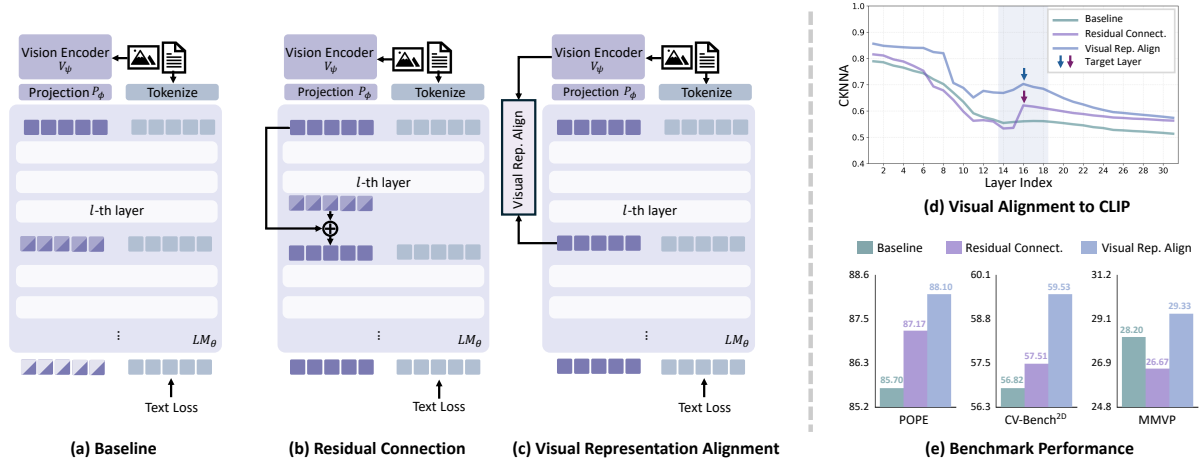


Figure 2: **Re-injecting or aligning visual features improves representation alignment and performance:** (a–c) Comparison of (a) baseline visual instruction tuning (Liu et al., 2023), (b) re-injecting visual features, and (c) visual representation alignment, all applied at the 16th layer. (d) Layer-wise alignment between visual tokens in MLLMs and vision encoder features, measured by CKNNA (Huh et al., 2024), with shaded regions denoting middle layers that are particularly important for visual understanding. (e) Benchmark performance corresponding to (a–c).

that can undermine performance on tasks requiring complex visual reasoning or grounding.

To empirically validate this hypothesis, we measure the similarity between the internal visual representations of LLaVA (Liu et al., 2024a) and the original visual features \mathbf{z} extracted by its vision encoder (e.g., CLIP (Radford et al., 2021)). We adopt CKNNA (Huh et al., 2024) as a metric to quantify representational similarity.

As shown in Fig. 2-(d), similarity to CLIP features drops sharply after the early layers and remains low in deeper layers, indicating that the model’s internal visual representations increasingly diverge from the encoder’s input features. This trend suggests that, without explicit visual supervision, the model has little incentive to preserve the encoder’s rich visual information.

Interestingly, despite the overall decline in alignment, the middle layers show a clear attenuation of this trend, with even slight increase, suggesting that the network implicitly benefits from retaining visual representations at these depths when generating visually grounded answers. This observation aligns with prior analyses of information flow in MLLMs (Zhang et al., 2025b; Kaduri et al., 2025) and is also confirmed by our later layer-wise ablations, which show that leveraging the middle layers for vision-centric tasks shows the largest gains.

- **Remark 1.** Internal visual representations in MLLMs progressively lose rich visual information originally provided by the input vision encoder.

4.2 Does preserving visual information help?

Having observed the mid-layer local increase in representation alignment, we ask whether *explicitly preserving* such visual information is beneficial. Let $\mathbf{e}_{\ell}^{\text{img}} \in \mathbb{R}^{N \times D}$ denote the visual representations at the ℓ -th layer of MLLMs. As a direct approach (Fig. 2-(b)), we re-inject the projected visual representation $P_{\phi}(\mathbf{z})$ into an intermediate layer of the language model via a residual path:

$$\mathbf{e}_{\ell,i}^{\text{img}} \leftarrow \mathbf{e}_{\ell,i}^{\text{img}} + P_{\phi}(\mathbf{z}_i). \quad (2)$$

To isolate the effect of visual information retention without introducing new supervision, the model is trained solely with the original text loss \mathcal{L}_{LM} . Unless otherwise stated, we set $\ell = 16$ in a 32-layer model LLaVA (Liu et al., 2024a), as fine-grained visual understanding emerges most prominently in middle layers, consistent with later layer-wise ablations (see Sec. 5.3).

As shown in Fig. 2-(d), adding the residual connection better preserves the alignment with the encoder’s visual features, as indicated by higher CKNNA similarity. Evaluated across standard benchmarks (Fig. 2-(e)), this approach shows general improvements over the baseline, supporting the hypothesis that retaining encoder-aligned visual information benefits downstream tasks. Although the residual connection provides general gains, concerns remain that the residual connection often only provides local impact on the specific layer of injection rather than boosting the entire visual pathway of the MLLMs to better preserve

the rich visual information originally provided by the input vision encoder.

4.3 Visual Representation Alignment

Representation alignment with encoder features.

Beyond residual connection, we further explore a more principled approach, which is to *explicitly* regularize intermediate visual representations to align with the encoder features, which enables the regularization of a wider range of the visual pathway; see Fig. 2-(c). Let \mathbf{z} denote the frozen encoder features from $V_\psi(\cdot)$ and $\mathbf{e}_\ell^{\text{img}} \in \mathbb{R}^{N \times D}$ the visual representations at the ℓ -th layer of the MLLM. We introduce a learnable projection $P_\pi(\cdot)$ to map $\mathbf{e}_\ell^{\text{img}}$ into the encoder feature space and define the visual representation alignment loss:

$$\mathcal{L}_{\text{VRA}}(\mathbf{e}_\ell^{\text{img}}, \mathbf{z}) = -\frac{1}{N} \sum_{i=1}^N \text{sim}\left(P_\pi(\mathbf{e}_{\ell,i}^{\text{img}}), \mathbf{z}_i\right), \quad (3)$$

where $\text{sim}(\cdot, \cdot)$ is cosine similarity and gradients do not flow into \mathbf{z} . Finally, the total objective augments the language modeling loss with this alignment term:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{LM}} + \lambda \mathcal{L}_{\text{VRA}}, \quad (4)$$

with λ controlling the strength of alignment.

As shown in Fig. 2-(d,e), this alignment outperforms residual connection in all multimodal benchmarks while also evidenced by the higher CKNNA similarity. Further analysis on this finding is provided in Appx. B. This shows that constraining intermediate features through an alignment loss offers stronger preservation of fine-grained semantics through explicit regularization, while residual connections offers only weak constraints without enforcing consistency at the feature level.

- **Remark 2.** Preventing visual information loss at the intermediate visual representation enhances visual understanding capabilities in MLLMs.

Despite the general performance boost from retaining encoder-aligned visual information, either by re-injecting projected features or applying visual representation alignment, a notable exception is MMVP (Tong et al., 2024b), which targets cases where CLIP-like features underperform. In this setting, performance shows only marginal improvement or even a slight drop, suggesting that propagating the encoder’s features can also transmit its

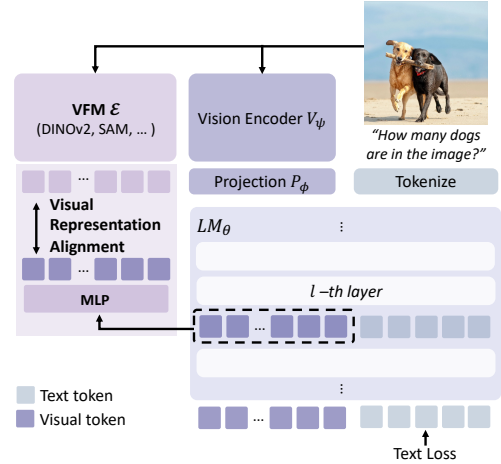


Figure 3: **Illustration of VIRAL.** We align visual pathway representation from MLLMs to strong, informative representations from VFMs to improve the vision understanding performance of MLLMs.

inductive biases and limitations. These findings raise the question of the *alignment target*: should the model remain tied to the original encoder features \mathbf{z} , or be guided toward more informative visual semantics? While aligning to \mathbf{z} helps retain meaningful attributes, its utility is constrained by the encoder’s representational capacity.

From encoder features to other VFMs. Motivated by this, we adopt stronger vision foundation models (VFMs) as teachers to supervise internal visual representations, providing richer vision-centric targets that complement language supervision. Building on this insight, we propose **Visual Representation Alignment (VIRAL)**, which aligns intermediate MLLM visual representations with features from a pretrained VFM, thereby preserving richer visual semantics than those available from the encoder alone. Let $\mathcal{E}(\cdot)$ denote a pretrained VFM encoder. Given an input image I , the encoder produces target features $\mathbf{y} = \mathcal{E}(I) \in \mathbb{R}^{N \times d}$, where d is the VFM feature dimension. Let $\mathbf{e}_\ell^{\text{img}} \in \mathbb{R}^{N \times D}$ be the MLLM’s visual representations at layer ℓ , and let $P_\pi(\cdot)$ be a learnable projection that maps $\mathbf{e}_\ell^{\text{img}}$ into the VFM feature space. We instantiate the visual representation alignment loss by replacing the encoder target \mathbf{z} in Eq. 3 with \mathbf{y} :

$$\mathcal{L}_{\text{VRA}}(\mathbf{e}_\ell^{\text{img}}, \mathbf{y}) = -\frac{1}{N} \sum_{i=1}^N \text{sim}\left(P_\pi(\mathbf{e}_{\ell,i}^{\text{img}}), \mathbf{y}_i\right). \quad (5)$$

Minimizing \mathcal{L}_{VRA} regularizes the MLLM’s internal visual pathway to align with the VFM. The

Language Model	Vision Encoder	\mathcal{L}_{VRA}	CV-Bench ^{2D}	MMVP	What's Up	Hallu. Bench	POPE	MMStar	MME	MMMU	MMBench ^{En}	MMVet	AI2D	OCRBench
Vicuna-1.5-7B	CLIP	✗	56.82	28.20	40.13	44.48	85.70	33.93	1650.21	33.78	63.23	26.65	53.82	319.0
		✓	59.67	33.33	48.55	46.06	88.32	33.93	1694.52	35.33	63.06	29.45	54.37	315.0
	SigLIPv2	✗	58.90	28.22	40.90	44.69	90.13	36.53	1738.96	34.00	64.00	29.31	55.76	384.0
		✓	62.66	33.11	44.40	45.95	90.77	37.20	1835.62	37.66	64.18	30.41	56.15	383.0
Qwen2.5-7B	CLIP	✗	58.97	33.47	59.08	44.37	85.88	39.20	1743.56	38.22	64.18	26.42	61.46	306.0
		✓	60.50	36.07	63.57	45.01	84.92	39.67	1765.65	41.44	68.13	28.62	64.12	310.0
Vicuna-1.5-13B	CLIP	✗	57.51	32.30	44.44	41.74	87.12	34.47	1599.04	34.89	66.07	31.47	57.67	334.0
		✓	58.97	37.80	62.26	42.06	87.79	37.00	1636.62	37.22	66.07	32.39	56.44	336.0

Table 1: **Effects of visual representation alignment.** We compare models trained with and without \mathcal{L}_{VRA} across various vision encoders and LLM backbones, evaluating them on both vision-centric and general multimodal benchmarks. Our simple regularization, \mathcal{L}_{VRA} , combined with DINOv2, consistently improves performance.

overall framework is illustrated in Fig. 3.

5 Experiments

5.1 Experimental Settings

Implementation details. We build on the widely used LLaVA-1.5 (Liu et al., 2024a), which provides fully open-source model weights and training data, enabling controlled and reproducible analysis. We used Vicuna-1.5 (Chiang et al., 2023) as the language model with a CLIP vision encoder (Radford et al., 2021). Following its instruction-tuning recipe, we adopt LoRA (Hu et al., 2022) for efficient adaptation as prior work reports that LLaVA-1.5 with LoRA attains comparable performance to full fine-tuning (Liu et al., 2024a). Unless otherwise noted, we use the original LLaVA-665K dataset without any additional data. The visual-representation projector $P_{\pi}(\cdot)$ is a lightweight three-layer MLP with SiLU activations, and we set $\mathcal{E}(\cdot)$ to DINOv2 as default (Sec. 5.3).

Evaluation. To demonstrate the effectiveness of VIRAL, we evaluate it on widely used benchmarks across four categories, which together assess whether our method improves vision-centric and hallucination-sensitive performance without compromising general MLLM capabilities: (1) vision-centric tasks requiring spatial reasoning or object counting, including CV-Bench^{2D} (Tong et al., 2024a), What’s Up (Chen et al., 2025; Kamath et al., 2023), and MMVP (Tong et al., 2024b); (2) multimodal hallucination detection, using POPE (Li et al., 2023) and Hallusion-bench (Guan et al., 2024); (3) general multimodal understanding, including MME (Yin et al., 2024), MMStar (Chen et al., 2024a), MMMU (Yue et al., 2024), MMBench (Liu et al., 2024c), MMVet (Yu

et al., 2023), as well as (4) document understanding benchmarks such as AI2D (Kembhavi et al., 2016), and OCRBench (Fu et al., 2024). Across all benchmarks, we adopt Imms-eval (Zhang et al., 2024) when available, and otherwise follow the original benchmark protocols. Detailed implementation details are provided in Appx. A.

5.2 Main Results

Tab. 1 summarizes results on vision-centric, hallucination, general vision–language benchmarks and document understanding benchmarks.

LLaVA-1.5-7B. With identical training settings, VIRAL consistently outperforms the baseline, yielding substantial improvements on fine-grained vision-centric tasks while preserving general multimodal performance, including document understanding, via intermediate feature alignment with VFM targets.

With stronger vision encoders. To demonstrate that our method goes beyond simply adopting a stronger vision encoder or input features, we evaluate our method using an architecture equipped with SigLIPv2 (Tschannen et al., 2025), trained with both contrastive (CLIP-style) and self-supervised (DINO-style) objectives. Even with this stronger encoder, our alignment loss yields consistent improvements, showing that the gains stem from the alignment itself.

Robustness to language backbones We additionally evaluate a scaled-up language backbone, comparing Vicuna-1.5-13B with its 7B counterpart, as well as an alternative backbone, Qwen2.5-7B (Bai et al., 2025), to demonstrate that our method is not tied to a specific language model.

VFM	CV-Bench ^{2D}	MMVP	What's Up	POPE	MME
Baseline	56.82	28.20	40.13	85.70	1650.21
DINOv2	59.67	33.33	48.55	88.32	1694.52
CLIP	59.53	29.33	44.50	88.10	1548.49
SAM	57.58	30.27	49.84	88.34	1648.77
DAv2	58.55	33.33	47.29	88.70	1682.42
RADIO	57.59	35.33	47.35	88.52	1692.94

Table 2: **Effects of different VFMs.** Performance of \mathcal{L}_{VRA} with different VFMs across benchmarks.

Layer	CV-Bench ^{2D}	MMVP	What's Up	POPE	MME
Baseline	56.82	28.20	40.13	85.70	1650.21
4	58.55	30.67	45.05	87.68	1720.36
8	58.28	27.70	48.32	88.43	1662.67
12	57.77	28.59	48.19	88.27	1648.88
16	59.67	33.33	48.55	88.32	1694.52
20	55.22	27.41	48.04	88.39	1705.97
24	55.77	27.48	47.99	88.10	1740.55
28	54.87	27.19	47.82	88.56	1755.86
32	56.12	26.52	47.60	87.32	1678.69

Table 3: **Effects of target layers.** Performance of \mathcal{L}_{VRA} with different target layers across benchmarks.

These results indicate that regularizing intermediate visual representations is a generally applicable strategy that improves MLLMs across vision encoders, model scales, and language backbones.

5.3 Component-wise Analysis

In this ablation study, we conduct a comprehensive analysis of key design choices underlying our framework, focusing on core components: the selection of target visual features and the choice of alignment layer. We evaluate the impact of each component across five benchmarks (CV-Bench, MMVP, What's Up, POPE, and MME) to validate their respective contributions to the model's performance on vision-grounded tasks. Unless otherwise specified, all baseline results correspond to LLaVA-1.5-7B. Additional ablation studies on alignment objectives and target layers are provided in Appx. C.

Vision foundation models. We begin by identifying the most effective target visual features for aligning internal visual representations in MLLMs, as summarized in Tab. 2. While residual connections and alignment with CLIP (LLaVA-1.5's original vision encoder) help improve visual comprehension (Fig. 2), their performance on spatial tasks like MMVP is limited—likely due to CLIP's weakness in modeling spatial relations (Yuksekonul et al., 2022). To address this, we evaluate several stronger vision foundation models (VFMs), including DINOv2 (Oquab et al., 2023), CLIP (Radford et al., 2021), Segment Anything (Kirillov et al., 2023) (SAM), Depth Anything v2 (Yang et al., 2024) (DAv2), and RADIOv2.5 (Heinrich et al.,

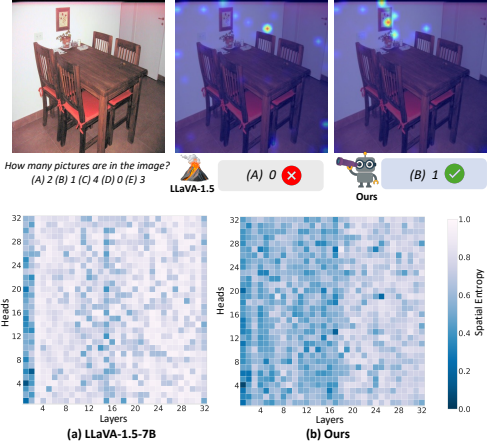


Figure 4: **Analysis of attention.** Qualitative comparison on text-to-image attention maps (top) and quantified spatial entropy of attention across layers and heads (bottom). Applying VIRAL encourages model to attend to more contextually important content, yielding a more focused and structured attention pattern.

2025). Our analysis confirms that aligning with stronger visual features further enhances visual understanding. Results show that DINOv2 consistently emerges as the most effective and versatile, and we thus adopt DINOv2 as the default visual foundation model for all experiments.

Target layers. We then analyze alignment at individual target layers to determine the most effective position as shown in Tab. 3. Here we report performance at every 4th layer throughout the network. We observe that performance varies depending on the alignment layer, with the 16th layer of the 32-layer model consistently yielding stronger results across multiple benchmarks. This trend is consistent with prior findings (Zhang et al., 2025b; Kaduri et al., 2025) and our earlier analysis, suggesting that certain layers in MLLMs are particularly attuned to visual information processing.

5.4 Attention Analysis

We analyze the effectiveness of our proposed framework with visual representation alignment in terms of text-to-image attention, as shown in Fig. 4 (top). The attention map produced by the \mathcal{L}_{VRA} trained model exhibits more semantically aligned focus on image regions corresponding to the given textual prompts. To quantify this, we adopt spatial entropy (Batty, 1974), motivated by (Kang et al., 2025b), as a metric of attention localization. As shown in Fig. 4 (bottom), LLaVA-1.5-7B exhibits high entropy across layers and heads, reflecting dispersed attention patterns, whereas our model

Vision Enc.	\mathcal{L}_{VRA}	original	patch shuffle	Δ
CLIP	✗	400	374	-26 (6.5%)
	✓	414	360	-54 (13.0%)
SigLIPv2	✗	374	353	-21 (5.6%)
	✓	436	353	-83 (19.0%)

Table 4: **Robustness to token permutation.** Number of correct predictions out of 788 spatial reasoning tasks in CV-Bench^{2D}.

shows consistently lower entropy—particularly at the aligned intermediate layer—indicating more selective and meaningful attention patterns.

5.5 Robustness Analysis

We investigate whether VIRAL enables MLLMs to better capture spatial relationships. Prior work (Qi et al., 2025) shows that many MLLMs are weakly grounded in 2D spatial structure: even when visual tokens are randomly permuted—effectively destroying image layout (see Appx. A)—performance drops only marginally, suggesting an order-insensitive, bag-of-patches representation. To assess whether our method increases sensitivity to spatial cues, we permute visual tokens $\mathbf{z} = V_\psi(I)$ before feeding them into $\text{LM}_\theta(\cdot)$ and evaluate performance on the spatial reasoning category of CV-Bench^{2D}. As shown in Tab. 4, while the text-only supervised baseline exhibits minimal degradation, our model suffers substantially larger drops, indicating stronger reliance on spatial structure. This confirms that our loss encourages MLLMs to better exploit fine-grained spatial relationships.

5.6 Qualitative Analysis

We qualitatively demonstrate the effectiveness of VIRAL through detailed analyses of model outputs and internal visual representations. By adopting VIRAL, we observe substantial improvements in performance on vision-centric tasks. As illustrated in Fig. 5, VIRAL correctly answers challenging visual questions related to the number of objects and spatial positioning, whereas the baseline model, LLaVA-1.5-7B, frequently fails.

Furthermore, by aligning internal visual representations with robust vision foundation models (VFMs), the semantic quality of intermediate representations is significantly enhanced. This improvement is clearly evidenced in the PCA visualizations from the visual representations obtained from the 16-th layer of Ours and LLaVA-1.5-7B shown in Fig. 5. These visualizations highlight that VIRAL effectively guides the model to preserve critical visual details, thereby facilitating better fine-grained

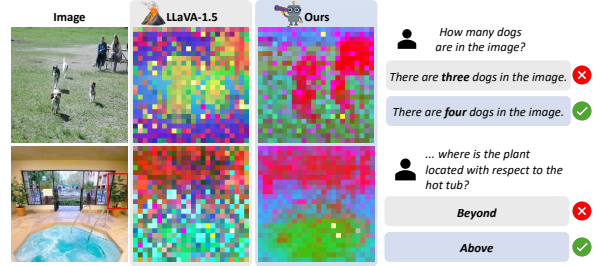


Figure 5: **Qualitative comparison of baseline and VIRAL.** The first column shows the input image–question pairs, and the next two present LLaVA-1.5 and VIRAL results with PCA visualizations and answers. VIRAL yields structured embeddings and correct answers on counting and spatial tasks where the baseline fails.

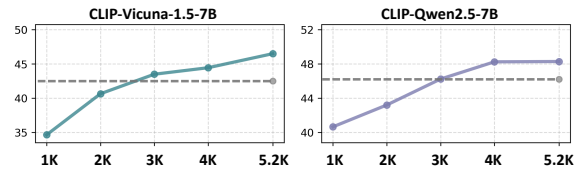


Figure 6: **Training Efficiency.** Performance with \mathcal{L}_{VRA} (solid) evaluated every 1K steps, averaging accuracies on CV-Bench^{2D} and MMVP. Models trained with \mathcal{L}_{VRA} achieve faster convergence. Dashed lines represent converged performance of baseline.

visual comprehension. Additional visualizations are provided in Appx. G.1 and G.2.

5.7 Training Efficiency

To further showcase the benefits of VIRAL, we evaluate vision-centric benchmarks, including CV-Bench^{2D} and MMVP, averaging accuracy every 1K training steps over the 5.2K steps of visual instruction tuning (Fig. 6). Both CLIP-based models trained with \mathcal{L}_{VRA} converge faster and surpass baseline performance within 3K steps. Since our method introduces only about a 3% overhead in total training time, these early gains can translate into improved scalability. This indicates that our approach enhances final accuracy while also accelerating training.

6 Conclusion

In this work, we propose VIRAL, a simple yet effective regularization strategy that aligns the internal visual representations of MLLMs with those from pre-trained vision foundation models. Our approach helps preserve fine-grained visual semantics often discarded under text-only supervision, thereby enabling more accurate spatial reasoning and object grounding.

582 Limitations

583 Our approach introduces an additional training-
584 time component by incorporating external VFMs
585 to regularize internal visual representations. While
586 this alignment is applied only during training and
587 does not affect inference cost, it may incur a mod-
588 est increase in computational overhead during op-
589 timization, particularly when high-capacity VFMs
590 are used. In addition, the benefits of our method
591 tend to be task-dependent. Since the alignment en-
592 courages visual representations to capture rich se-
593 mantic and structural cues distilled from the VFM,
594 performance gains may be limited on tasks whose
595 primary signals are weakly correlated with such
596 semantics, such as OCR-centric benchmarks. Nev-
597 ertheless, we do not observe performance degrada-
598 tion on these tasks; the improvements instead tend
599 to be marginal compared to those on vision-centric
600 benchmarks.

601 Ethical considerations

602 This work introduces a training-time regulariza-
603 tion method for aligning visual representations in
604 multimodal large language models. It does not in-
605 volve new data collection or user interaction and
606 therefore raises no additional privacy or security
607 concerns beyond those inherent to existing multi-
608 modal models.

609 References

610 Mozghan Nasr Azadani, James Riddell, Sean Sedwards,
611 and Krzysztof Czarnecki. 2025. Leo: Boosting mix-
612 ture of vision encoders for multimodal large language
613 models. *arXiv preprint arXiv:2501.06986*.

614 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,
615 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei
616 Huang, and 1 others. 2023a. Qwen technical report.
617 *arXiv preprint arXiv:2309.16609*.

618 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang,
619 Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou,
620 and Jingren Zhou. 2023b. [Qwen-vl: A versa-
621 tile vision-language model for understanding, lo-
622 calization, text reading, and beyond](#). *Preprint*,
623 [arXiv:2308.12966](#).

624 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-
625 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie
626 Wang, Jun Tang, and 1 others. 2025. Qwen2. 5-vl
627 technical report. *arXiv preprint arXiv:2502.13923*.

628 Michael Batty. 1974. Spatial entropy. *Geographical
629 analysis*, 6(1):1–31.

Daniel Bolya, Po-Yao Huang, Peize Sun, Jang Hyun
630 Cho, Andrea Madotto, Chen Wei, Tengyu Ma, Jiale
631 Zhi, Jathushan Rajasegaran, Hanoona Rasheed, and
632 1 others. 2025. Perception encoder: The best vi-
633 sual embeddings are not at the output of the network.
634 *arXiv preprint arXiv:2504.13181*. 635

Junbum Cha, Wooyoung Kang, Jonghwan Mun, and
636 Byungseok Roh. 2024. Honeybee: Locality-
637 enhanced projector for multimodal llm. In *Proceeed-
638 ings of the IEEE/CVF Conference on Computer Vi-
639 sion and Pattern Recognition*, pages 13817–13827. 640

Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang
641 Zang, Zehui Chen, Haodong Duan, Jiaqi Wang,
642 Yu Qiao, Dahua Lin, and 1 others. 2024a. Are we
643 on the right way for evaluating large vision-language
644 models? *Advances in Neural Information Processing
645 Systems*, 37:27056–27087. 646

Shiqi Chen, Tongyao Zhu, Ruochen Zhou, Jinghan
647 Zhang, Siyang Gao, Juan Carlos Niebles, Mor Geva,
648 Junxian He, Jiajun Wu, and Manling Li. 2025. Why
649 is spatial reasoning hard for vlms? an attention mech-
650 anism perspective on focus areas. *arXiv preprint
651 arXiv:2503.01773*. 652

Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu,
653 Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong
654 Ye, Hao Tian, Zhaoyang Liu, and 1 others. 2024b.
655 Expanding performance boundaries of open-source
656 multimodal models with model, data, and test-time
657 scaling. *arXiv preprint arXiv:2412.05271*. 658

Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo
659 Chen, Sen Xing, Muyan Zhong, Qinglong Zhang,
660 Xizhou Zhu, Lewei Lu, and 1 others. 2024c. Internvl:
661 Scaling up vision foundation models and aligning
662 for generic visual-linguistic tasks. In *Proceedings of
663 the IEEE/CVF conference on Computer Vision and
664 Pattern Recognition*, pages 24185–24198. 665

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,
666 Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan
667 Zhuang, Yonghao Zhuang, Joseph E. Gonzalez,
668 Ion Stoica, and Eric P. Xing. 2023. Vicuna:
669 An open-source chatbot impressing gpt-4 with
670 90% chatgpt quality. [https://lmsys.org/blog/
671 2023-03-30-vicuna/](https://lmsys.org/blog/2023-03-30-vicuna/). Accessed: 2025-08-19. 672

Ling Fu, Zhebin Kuang, Jiajun Song, Mingxin Huang,
673 Biao Yang, Yuzhe Li, Linghao Zhu, Qidi Luo, Xinyu
674 Wang, Hao Lu, and 1 others. 2024. Ocrbench v2:
675 An improved benchmark for evaluating large multi-
676 modal models on visual text localization and reason-
677 ing. *arXiv preprint arXiv:2501.00321*. 678

Stephanie Fu, Tyler Bonnen, Devin Guillory, and Trevor
679 Darrell. 2025. Hidden in plain sight: Vlms over-
680 look their visual representations. *arXiv preprint
681 arXiv:2506.08008*. 682

Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian,
683 Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen,
684 Furong Huang, Yaser Yacoob, and 1 others. 2024.
685 Hallusionbench: an advanced diagnostic suite for
686

687	entangled language hallucination and visual illusion in large vision-language models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 14375–14385.	Tombari. 2024. Brave: Broadening the visual encoding of vision-language models. In <i>European Conference on Computer Vision</i> , pages 113–132. Springer.	741
688			742
689			743
690			
691	Greg Heinrich, Mike Ranzinger, Hongxu Yin, Yao Lu, Jan Kautz, Andrew Tao, Bryan Catanzaro, and Pavlo Molchanov. 2025. Radiov2. 5: Improved baselines for agglomerative vision foundation models. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 22487–22497.	Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. 2016. A diagram is worth a dozen images. In <i>European conference on computer vision</i> , pages 235–251. Springer.	744
692			745
693			746
694			747
695			748
696			
697	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. <i>ICLR</i> , 1(2):3.	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, and 1 others. 2023. Segment anything. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 4015–4026.	749
698			750
699			751
700			752
			753
			754
701	Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In <i>Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition</i> , pages 6700–6709.	Black Forest Labs. 2025. The learnability-quality-compression trade-off. https://bfl.ai/techblog/representation-comparison/#ref-dieleman2025latents .	755
702			756
703			757
704			758
705			
706	Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. 2024. The platonic representation hypothesis. <i>arXiv preprint arXiv:2405.07987</i> .	Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and 1 others. 2024. Llava-onevision: Easy visual task transfer. <i>arXiv preprint arXiv:2408.03326</i> .	759
707			760
708			761
			762
			763
709	Jitesh Jain, Zhengyuan Yang, Humphrey Shi, Jianfeng Gao, and Jianwei Yang. 2025. Elevating visual perception in multimodal llms with visual embedding distillation. In <i>The Thirty-ninth Annual Conference on Neural Information Processing Systems</i> .	Wenyan Li, Raphael Tang, Chengzu Li, Caiqi Zhang, Ivan Vulić, and Anders Søgaard. 2025. Lost in embeddings: Information loss in vision-language models. <i>Preprint</i> , arXiv:2509.11986.	764
710			765
711			766
712			767
713			
714	Zhangqi Jiang, Junkai Chen, Beier Zhu, Tingjin Luo, Yankun Shen, and Xu Yang. 2025. Devils in middle layers of large vision-language models: Interpreting, detecting and mitigating object hallucinations via attention lens. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 25004–25014.	Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Evaluating object hallucination in large vision-language models. <i>arXiv preprint arXiv:2305.10355</i> .	768
715			769
716			770
717			771
718			
719			772
720			773
			774
			775
			776
721	Omri Kaduri, Shai Bagon, and Tali Dekel. 2025. What’s in the image? a deep-dive into the vision of vision language models. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 14549–14558.	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In <i>Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition</i> , pages 26296–26306.	777
722			778
723			779
724			
725			
726	Amita Kamath, Jack Hessel, and Kai-Wei Chang. 2023. What’s “up” with vision-language models? investigating their struggle with spatial reasoning. <i>arXiv preprint arXiv:2310.19785</i> .	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual instruction tuning. <i>Advances in neural information processing systems</i> , 36:34892–34916.	780
727			781
728			782
729			783
730	Seil Kang, Jinyeong Kim, Junhyeok Kim, and Seong Jae Hwang. 2025a. See what you are told: Visual attention sink in large multimodal models. <i>arXiv preprint arXiv:2503.03321</i> .	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, and 1 others. 2024c. Mmbench: Is your multi-modal model an all-around player? In <i>European conference on computer vision</i> , pages 216–233. Springer.	784
731			785
732			786
733			787
			788
			789
734	Seil Kang, Jinyeong Kim, Junhyeok Kim, and Seong Jae Hwang. 2025b. Your large vision-language model only needs a few attention heads for visual grounding. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 9339–9350.	Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, and 1 others. 2024. Deepseek-vl: towards real-world vision-language understanding. <i>arXiv preprint arXiv:2403.05525</i> .	790
735			791
736			792
737			793
738			794
739	Oğuzhan Fatih Kar, Alessio Tonioni, Petra Poklucar, Achin Kulshrestha, Amir Zamir, and Federico		
740			

795	Zixian Ma, Jerry Hong, Mustafa Omer Gul, Mona Gandhi, Irena Gao, and Ranjay Krishna. 2023. Crepe: Can vision-language foundation models reason compositionally? In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 10910–10921.	Peter Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Adithya Jairam Vedagiri IYER, Sai Charitha Akula, Shusheng Yang, Jihan Yang, Manoj Middepogu, Ziteng Wang, and 1 others. 2024a. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. <i>Advances in Neural Information Processing Systems</i> , 37:87310–87356.	850 851 852 853 854 855 856
801	Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruvi Shah, Xianzhi Du, Futang Peng, Anton Belyi, and 1 others. 2024. Mm1: methods, analysis and insights from multimodal llm pre-training. In <i>European Conference on Computer Vision</i> , pages 304–323. Springer.	Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. 2024b. Eyes wide shut? exploring the visual shortcomings of multimodal llms. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 9568–9578.	857 858 859 860 861 862
808	Clement Neo, Luke Ong, Philip Torr, Mor Geva, David Krueger, and Fazl Barez. 2024. Towards interpreting visual information processing in vision-language models. <i>arXiv preprint arXiv:2410.07149</i> .	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	863 864 865 866 867 868
812	OpenAI. 2023. Gpt-4v(ision) technical work and authors. https://openai.com/contributions/gpt-4v/ . Accessed: 2025-08-02.	Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, and 1 others. 2025. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. <i>arXiv preprint arXiv:2502.14786</i> .	869 870 871 872 873 874 875
815	Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, and 1 others. 2023. Dinov2: Learning robust visual features without supervision. <i>arXiv preprint arXiv:2304.07193</i> .	Pavan Kumar Anasosalu Vasu, Fartash Faghri, Chun-Liang Li, Cem Koc, Nate True, Albert Antony, Gokula Santhanam, James Gabriel, Peter Grusch, Oncel Tuzel, and 1 others. 2025. Fastvlm: Efficient vision encoding for vision language models. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 19769–19780.	876 877 878 879 880 881 882
821	Jianing Qi, Jiawei Liu, Hao Tang, and Zhigang Zhu. 2025. Beyond semantics: Rediscovering spatial awareness in vision-language models. <i>arXiv preprint arXiv:2503.17349</i> .	Constantin Venhoff, Ashkan Khakzar, Sonia Joseph, Philip Torr, and Neel Nanda. 2025. How visual representations map to language feature space in multimodal llms. <i>arXiv preprint arXiv:2506.11976</i> .	883 884 885 886
825	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, and 1 others. 2021. Learning transferable visual models from natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PmLR.	Gaurav Verma, Minje Choi, Kartik Sharma, Jamelle Watson-Daniels, Sejoon Oh, and Srijan Kumar. 2024. Cross-modal projection in multimodal llms doesn't really project visual attributes to textual space. <i>arXiv preprint arXiv:2402.16832</i> .	887 888 889 890 891
832	Mike Ranzinger, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. 2024. Am-radio: Agglomerative vision foundation model reduce all domains into one. In <i>Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition</i> , pages 12490–12500.	Haochen Wang, Anlin Zheng, Yucheng Zhao, Tiancai Wang, Zheng Ge, Xiangyu Zhang, and Zhaoxiang Zhang. 2024. Reconstructive visual instruction tuning. <i>arXiv preprint arXiv:2410.09575</i> .	892 893 894 895
838	Min Shi, Fuxiao Liu, Shihao Wang, Shijia Liao, Subhashree Radhakrishnan, Yilin Zhao, De-An Huang, Hongxu Yin, Karan Sapra, Yaser Yacoob, and 1 others. 2024. Eagle: Exploring the design space for multimodal llms with mixture of encoders. <i>arXiv preprint arXiv:2408.15998</i> .	Zichen Wen, Yifeng Gao, Shaobo Wang, Junyuan Zhang, Qintong Zhang, Weijia Li, Conghui He, and Linfeng Zhang. 2025. Stop looking for important tokens in multimodal language models: Duplication matters more. <i>arXiv preprint arXiv:2502.11494</i> .	896 897 898 899 900
844	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, and 1 others. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. 2024. Depth anything v2. <i>Advances in Neural Information Processing Systems</i> , 37:21875–21911.	901 902 903 904

905 Senqiao Yang, Yukang Chen, Zhuotao Tian, Chengyao
906 Wang, Jingyao Li, Bei Yu, and Jiaya Jia. 2025. Vi-
907 sionzip: Longer is better but not necessary in vision
908 language models. In *Proceedings of the Computer
909 Vision and Pattern Recognition Conference*, pages
910 19792–19802.

911 Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing
912 Sun, Tong Xu, and Enhong Chen. 2024. A survey on
913 multimodal large language models. *National Science
914 Review*, 11(12):nwae403.

915 Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon
916 Jeong, Jonathan Huang, Jinwoo Shin, and Saining
917 Xie. 2024. Representation alignment for generation:
918 Training diffusion transformers is easier than you
919 think. In *The Thirteenth International Conference on
920 Learning Representations*.

921 Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang,
922 Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan
923 Wang. 2023. Mm-vet: Evaluating large multimodal
924 models for integrated capabilities. *arXiv preprint
925 arXiv:2308.02490*.

926 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng,
927 Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang,
928 Weiming Ren, Yuxuan Sun, and 1 others. 2024.
929 Mmmu: A massive multi-discipline multimodal un-
930 derstanding and reasoning benchmark for expert agi.
931 In *Proceedings of the IEEE/CVF Conference on Com-
932 puter Vision and Pattern Recognition*, pages 9556–
933 9567.

934 Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri,
935 Dan Jurafsky, and James Zou. 2022. When and
936 why vision-language models behave like bags-of-
937 words, and what to do about it? *arXiv preprint
938 arXiv:2210.01936*.

939 Junyi Zhang, Charles Herrmann, Junhwa Hur, Luisa
940 Polania Cabrera, Varun Jampani, Deqing Sun, and
941 Ming-Hsuan Yang. 2023. A tale of two features: Sta-
942 ble diffusion complements dino for zero-shot seman-
943 tic correspondence. *Advances in Neural Information
944 Processing Systems*, 36:45533–45547.

945 Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu,
946 Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuan-
947 han Zhang, Jingkang Yang, Chunyuan Li, and Zi-
948 wei Liu. 2024. *Lmms-eval: Reality check on the
949 evaluation of large multimodal models*. *Preprint*,
950 arXiv:2407.12772.

951 Xiangdong Zhang, Jiaqi Liao, Shaofeng Zhang, Fanqing
952 Meng, Xiangpeng Wan, Junchi Yan, and Yu Cheng.
953 2025a. Videorepa: Learning physics for video gen-
954 eration through relational alignment with foundation
955 models. *arXiv preprint arXiv:2505.23656*.

956 Zhi Zhang, Srishti Yadav, Fengze Han, and Ekaterina
957 Shutova. 2025b. Cross-modal information flow in
958 multimodal large language models. In *Proceedings
959 of the Computer Vision and Pattern Recognition Con-
960 ference*, pages 19781–19791.

Appendix

A Additional implementation details

All experiments in this paper are conducted on four NVIDIA A100 GPUs (40 GB each).

Vision foundation models. We use a diverse set of pretrained VFMs to supervise internal visual representations. DINOv2 (Oquab et al., 2023), CLIP (Radford et al., 2021), and Depth Anything v2 (Yang et al., 2024) (DAv2) are used as patch size 14 models, while RADIO-v2.5 (Heinrich et al., 2025) and SAM (Kirillov et al., 2023) are used as patch size 16 models. To match the 576 visual tokens produced by CLIP-ViT-L/14 at 336×336 resolution in LLaVA-1.5 (Liu et al., 2024a), we adopt the same resolution for patch size 14 models and resize inputs to 384×384 for patch size 16 models. For SAM, which expects 1024×1024 inputs, we pad the interpolated features to 1024×1024 and crop them to the region corresponding to the original image, following AM-RADIO (Ranzinger et al., 2024) to avoid quality degradation.

Loss function and weighting. The cosine similarity $\text{sim}(\mathbf{x}, \mathbf{y})$, as done in previous works, is computed as following $\text{sim}(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^\top \mathbf{y}}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2}$. To balance the alignment loss \mathcal{L}_{VRA} with the language modeling loss \mathcal{L}_{LM} , we set $\lambda = 0.5$ by default.

Benchmark settings. To demonstrate the effectiveness of VIRAL, we evaluate it on a broad set of widely used benchmarks, including CV-Bench, MMVP, What’s Up, MMStar, MME, MMMU, MMBench, MMVet, HallusionBench, POPE, AI2D, InfoVQA, and OCRBench. We use only the 2D subset of CV-Bench, as 3D tasks are beyond the scope of this work, and report overall accuracy on CV-Bench^{2D} rather than separately averaging ADE20K and COCO. For MMVP, we follow the standard evaluation protocol based on pair accuracy and report the average accuracy over 10 runs for stability. For POPE, we evaluate on COCO following LLaVA and report the average accuracy across the random and popular subsets. For What’s Up, we report the average accuracy over the COCO_{one} and COCO_{two} splits. For MME, we report MME^{EN} along with the summed scores of the perception and cognition categories. We use the en_dev split for MMBench, and the val splits for MMMU and InfoVQA. For all remaining benchmarks, we follow the default evaluation protocols and scoring metrics.

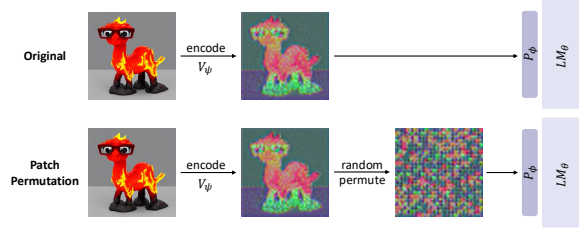


Figure 7: Visualization of patch random permutation experiments.

Spatial entropy. For Fig. 4, we compute average spatial entropy over generated text tokens. We use question–answer pairs from (Zhang et al., 2025b), which augment GQA (Hudson and Manning, 2019) with diverse categories and constrain answers to a single word or phrase. Among these, we focus on the Relation category and report the average spatial entropy within this subset.

Patch permutation. For our patch permutation experiment, we adopt the analysis pipeline originally proposed in (Qi et al., 2025). Specifically, we begin by extracting image features z from the vision encoder using $z = V_\psi(I)$, where I is the input image. Here, $z \in \mathbb{R}^{N \times H}$, with N denoting the number of visual tokens and H the dimensionality of the vision encoder features. Before processing the vision features z with the vision-language projector $P_\phi(\cdot)$ and language model $LM_\theta(\cdot)$, we apply a random permutation on the order of the visual tokens N , which is shown in the visualization of Fig. 7. This makes it extremely difficult to understand the visual attributes of the image, enabling us to evaluate how much the MLLM was understanding and utilizing the visual attributes originally available in the image.

B Extended exploration of the pilot study

In Sec. 4, we demonstrated that MLLMs exhibit progressive visual information loss across layers, and that preserving such information can enhance their visual understanding (Fig. 8-(b,c)). In this section, we compare two additional strategies for preserving visual information: a residual connection with the raw encoder feature prior to projection (*pre*-projection) as a direct approach to feature re-injection (Fig. 8-(d)), and our proposed visual representation alignment with the projected features (*post*-projection) provided to the language model (Fig. 8-(e)).

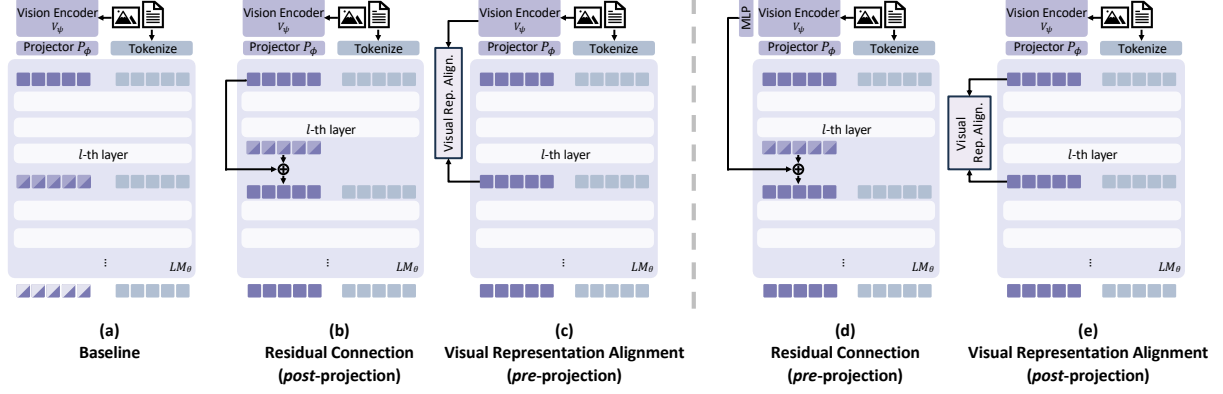


Figure 8: Extended exploration of the pilot study.

	POPE	CV-Bench ^{2D}	MMVP
Baseline	85.70	56.82	28.20
(b)	87.17	57.51	26.67
(c)	88.10	59.53	29.33
(d)	85.47	53.62	19.33
(e)	86.99	57.23	28.53

Table 5: Benchmark performance of the pilot study.

Residual connection with *pre*-projection features. Our investigation leverages *pre*-projection features—raw encoder features \mathbf{z} prior to the projector—through a direct residual connection to mitigate visual information loss within the language model, where a lightweight adapter $P_{\phi'}(\cdot)$ is employed for dimensional compatibility. As illustrated in Fig. 8-(d), we conduct one such experiment that re-injects \mathbf{z}_i into $\mathbf{e}_{\ell,i}^{\text{img}}$ such that

$$\mathbf{e}_{\ell,i}^{\text{img}} \leftarrow \mathbf{e}_{\ell,i}^{\text{img}} + P_{\phi'}(\mathbf{z}_i). \quad (6)$$

However, as shown in Tab. 5-(d), this approach generally performs worse than the baseline. This is because the raw encoder features, which have not passed through the pre-trained projector, are not sufficiently aligned with language features (Liu et al., 2023), and their direct residual connection consequently disrupts vision–language alignment in the intermediate layers. These findings suggest that incorporating external features into the internal visual pathway of LLMs requires more careful design.

Visual representation alignment with *post*-projection features. Next we further explore aligning the intermediate visual representation with the *post*-projection features, as shown in Fig. 8-(e). Here, we follow the same experimental setting as

in Sec. 4.3, while \mathcal{L}_{VRA} is defined as:

$$\mathcal{L}_{\text{VRA}} = -\frac{1}{N} \sum_{i=1}^N \text{sim}\left(P_{\pi}(\mathbf{e}_{\ell,i}^{\text{img}}), P_{\phi}(\mathbf{z}_i)\right). \quad (7)$$

The results presented in Tab. 5-(e) indicate that this approach generally improves performance over the baseline on vision-centric benchmarks, yet underperforms compared to leveraging raw features from the vision encoder. This may be attributed to the insufficient preservation of visual information in the *post*-projection features compared to the raw encoder outputs (Verma et al., 2024; Cha et al., 2024; Li et al., 2025).

C Additional ablation studies

Number of target layers. To investigate the effective number of target layers, we evaluate *multi-layer targets* around the 16th—specifically ± 1 (15–17) and ± 2 (14–18) ranges—and observe that applying alignment solely at the 16th layer achieves the best performance. These findings highlight that aligning visual representations at a specific pathway responsible for visual representation processing, rather than uniformly across multiple layers, is more effective in enhancing the visual understanding capabilities of MLLMs. Based on this observation, we adopt the 16th layer as the default alignment target with DINOv2.

Alignment objectives. We investigate the impact of different feature *alignment objectives* during instruction tuning. Specifically, we compare the performance of models trained with a feature relation alignment objective, as a substitute for the proposed direct visual representation alignment loss. Here, the alignment objective is defined as a

VFM	Layer Index	Objective	CV-Bench ^{2D}	MMVP	What’s Up	POPE	MME
Baseline			56.82	28.20	40.13	85.70	1650.21
<i>Ablation studies on different multi-layer targets</i>							
DINOv2	16	Cos. Sim.	59.67	33.33	48.55	88.32	1694.52
DINOv2	15 – 17	Cos. Sim.	59.32	28.00	47.17	87.61	1639.72
DINOv2	14 – 18	Cos. Sim.	49.62	22.55	42.58	87.90	1444.32
<i>Ablation studies on different alignment objectives</i>							
DINOv2	16	Cos. Sim.	59.67	33.33	48.55	88.32	1694.52
DINOv2	16	Relation	58.83	26.60	49.05	87.58	1674.30

Table 6: Ablation study on key design components.

mean squared error (MSE) loss between the self-similarity matrices of the VFM features and the transformed intermediate representations, which effectively distills the structural relationships among visual features following recent approaches (Zhang et al., 2025a; Bolya et al., 2025). As shown in Tab. 6, we find that simple cosine similarity-based alignment loss yields higher performance, and adopt it as our default strategy for alignment.

D Comparison with other training objectives

We compare our method with ROSS (Wang et al., 2024), which applies a reconstructive objective to the final hidden state of the visual representations. To isolate the sources of improvement, we implement two variants under identical experimental conditions: ROSS (Default), reproducing the original method, and ROSS (Middle), which applies the same objective to an intermediate layer (16th layer as in our configuration for target supervision).

Tab. 7 reveals several key findings that validate our approach. First, the critical importance of intermediate layer supervision—a contribution of our work—is evidenced by ROSS (Middle) outperforming ROSS (Default), particularly on vision-centric benchmarks. Although both ROSS (Default) and ROSS (Middle) show improvements over the baseline which also shows the importance of providing supervision to the visual pathways, the superiority of ROSS (Middle) over ROSS (Defaults) confirms our hypothesis that supervising visual information flow at strategically chosen intermediate layers, rather than naively at the model’s output, yields superior performance gains.

Second, and more fundamentally, our method significantly outperforms both ROSS variants across all benchmarks. This performance gap stems from a crucial distinction in objectives: while ROSS employs reconstruction-based objectives

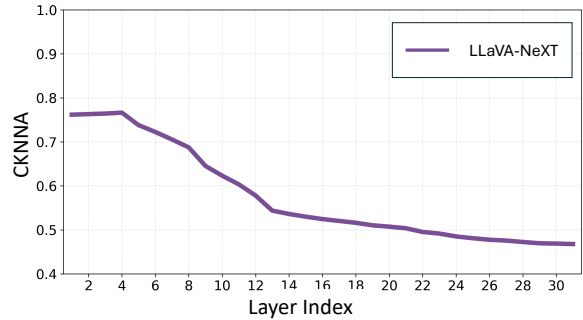


Figure 9: Visual alignment of LLaVA-NeXT (Liu et al., 2024b). Layer-wise alignment between visual tokens in MLLM and vision encoder features, measured by CKNNA and averaged across representations from tiled image splits.

that excel at preserving low-level fidelity, such approaches are inherently less suited for capturing the higher-level semantic abstractions required by complex reasoning tasks (Zhang et al., 2023; Tong et al., 2024a). In contrast, our direct alignment with pretrained vision foundation models provides richer semantic supervision that better bridges the vision-language gap.

These results demonstrate that our method’s superiority arises from two synergistic contributions: (1) the strategic placement of supervision at critical intermediate layers, and (2) the use of semantically-rich alignment signals from vision foundation models rather than reconstruction-based objectives. Together, these design choices enable more effective visual representation learning for multimodal understanding.

E Applicability of VIRAL to other MLLM architectures

Recent MLLMs employ various strategies to handle inputs with dynamic resolutions, including dynamically adjusting the sequence length of visual tokens (Bai et al., 2023b) or dividing high-resolution images into independently encoded tile

Language Model	Vision Encoder	Objective	CV-Bench ^{2D}	MMVP	What’s Up	POPE	MMStar	MME
Vicuna-1.5-7B	CLIP	Baseline	56.82	28.20	40.13	85.70	33.93	1650.21
		ROSS (Default)	54.24	29.73	43.57	88.19	34.73	1648.87
		ROSS (Middle)	56.05	31.40	45.98	88.21	33.53	1647.27
		VIRAL	59.67	33.33	48.55	88.32	33.93	1694.52

Table 7: Comparison with reconstructive objective.

Layer	CV-Bench ^{2D}	MMVP	What’s Up	POPE	MME
Baseline	56.82	28.20	40.13	86.90	1650
4	58.55	30.67	45.05	87.68	1720
8	58.28	27.70	48.32	88.43	1663
12	57.77	28.59	48.19	88.27	1649
16	59.67	33.33	48.55	88.32	1695
20	55.22	27.41	48.04	88.39	1706
24	55.77	27.48	47.99	88.10	1741
28	54.87	27.19	47.82	88.56	1756
32	56.12	26.52	47.60	87.32	1679

(a) DINOv2

Layer	CV-Bench ^{2D}	MMVP	What’s Up	POPE	MME
Baseline	56.82	28.20	40.13	86.90	1650
4	52.99	26.00	46.18	87.78	1614
8	57.58	24.66	44.42	87.88	1710
12	54.87	26.00	45.99	88.39	1672
16	57.58	30.27	49.84	88.34	1665
20	54.24	26.00	47.39	88.42	1698
24	57.79	26.66	48.87	87.79	1706
28	57.16	24.00	45.66	88.35	1698
32	57.09	29.33	44.42	87.80	1717

(b) SAM

Table 8: Layer selection ablation for different visual foundation models. All results are preliminary.

grids (Liu et al., 2024b; Chen et al., 2024b). The latter approach preserves the original image resolution and is commonly adopted to capture fine-grained visual details.

To investigate whether VIRAL can be applied to such recent MLLM paradigms, we examine if our core motivation—mitigating vision information loss—remains relevant within this tiled image processing strategy. Fig. 9 demonstrates a decline in alignment scores between input visual features and layer-wise visual representations across the layers in LLaVA-NeXT (Liu et al., 2024b), as measured using CKNNA (Huh et al., 2024). This shows similar patterns observed in Fig. 2(d), suggesting that VIRAL can be applied orthogonally to such techniques and has the potential to similarly enhance fine-grained visual understanding in MLLMs designed for dynamic resolution handling.

F Exploring target layers across different VFMs

Tab. 8 presents the target-layer ablation results when applying the \mathcal{L}_{VRA} with different VFMs, specifically DINOv2 and SAM. As shown in Tab. 2 we demonstrate that visual representation alignment not only with DINOv2 but also with alternative VFMs consistently outperforms the baseline, indicating that our approach is not restricted to a specific visual foundation model and can flexibly leverage diverse VFM features. In Tab. 2 and 6, we further show that selectively applying the \mathcal{L}_{VRA} to specific LLM layers is crucial for achieving effective alignment. Building on these findings, here we explore how the choice of the target layer can be specified across different VFMs, using DINOv2 and SAM as a representative examples. Across VFMs, we observe similar trends in target-layer selection. In particular, the layer ablation results on SAM also show that applying \mathcal{L}_{VRA} to middle layers generally leads to larger performance improvements. This suggests that, even with a broader choice of target VFMs, aligning visual representations at intermediate LLM layers, where visual information is more actively processed, can be effective.

G Additional visualizations and results

G.1 Layer-wise internal representations

We present PCA visualizations of the intermediate visual representations from all layers of LLaVA-1.5-7B and VIRAL in Fig. 10, enabling a layer-wise comparison of their representational structures. A qualitative comparison with the baseline reveals that visual representation alignment regularizes the MLLM’s internal visual features, leading to more semantically coherent and structured representation, especially in the middle and later layers where meaningful vision understanding emerges.

We present PCA visualizations of the intermediate visual representations from all layers of LLaVA-

1.5-7B and VIRAL in Fig. 10, enabling a layer-wise comparison of their representational structures. A qualitative comparison with the baseline reveals that visual representation alignment regularizes the MLLM’s internal visual features, leading to more semantically coherent and structured representation, especially in the middle and later layers where meaningful vision understanding emerges.

G.2 Visual representations with different VFMs

In addition to Fig. 5, we qualitatively present in Fig. 11 PCA visualizations of how internal visual representations evolve when aligned with different VFMs. Compared to the baseline representation from LLaVA-1.5-7B, VFM features exhibit more semantically structured organization. Aligning the MLLM’s internal representations with these VFM features distills such structure, enabling the model to refer to enhanced and more coherent visual representations.

G.3 Attention map visualizations

In Fig. 12, we provide visualizations of text-to-image cross-attention maps in the MLLM to qualitatively support the attention analysis from the main paper. Compared to the baseline, the model trained with our method exhibits improved attention behavior by focusing more accurately and locally on regions relevant to the given multimodal context. This observation aligns well with the spatial entropy analysis in Fig. 4, where models trained with visual representation alignment show more focused and discriminative attention patterns.

H Use of Large Language Models

We disclose that Large Language Models were used to assist in grammar correction and polishing of the writing in this paper.

I Potential risk

We do not introduce any additional potential risk beyond those inherent to the underlying architectures, but the integration of LLMs and VLMs involves potential risks concerning data provenance, algorithmic bias, and stochastic hallucinations. These models remain susceptible to adversarial exploitation—specifically prompt injection and input perturbations—which may circumvent safety alignment to generate erroneous or prejudicial outputs.



How many pillows are in the image?

(A) 2 (B) 1 (C) 4 (D) 0 (E) 3



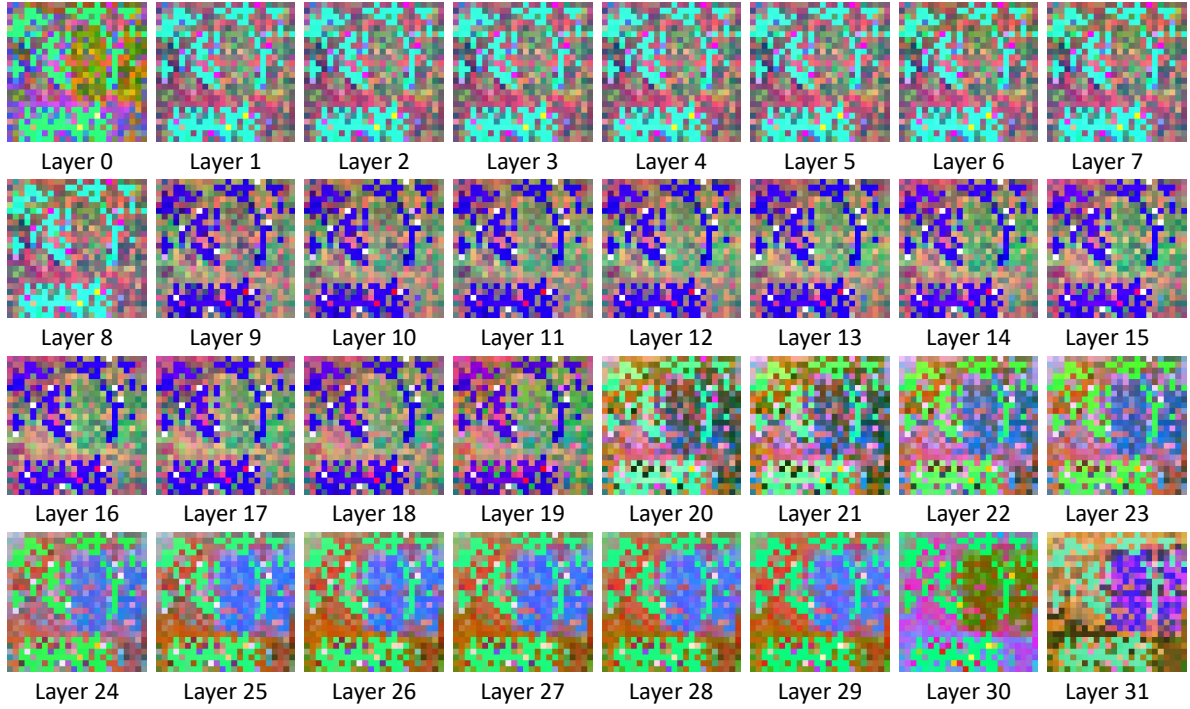
LLaVA-1.5

(D) 0

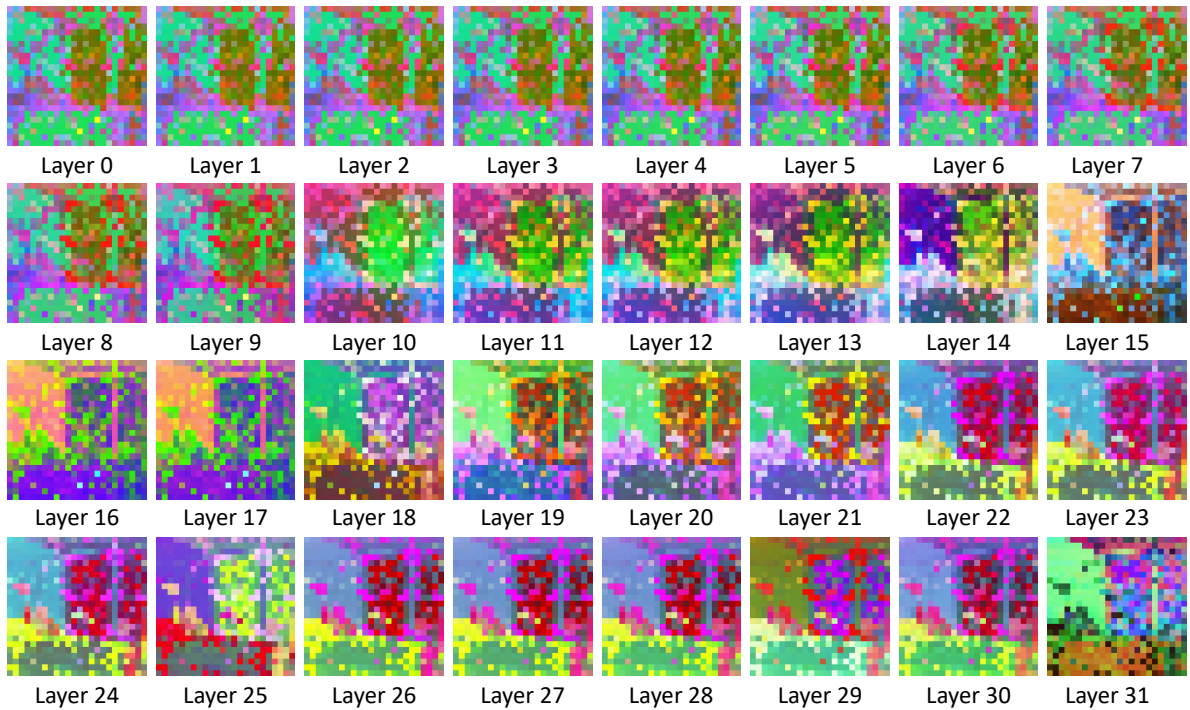


VIRAL

(A) 2



(a) LLaVA-1.5-7B



(b) LLaVA-1.5-7B + VIRAL

Figure 10: Layer-wise PCA visualizations of visual representations from (a) LLaVA-1.5-7B and (b) VIRAL.

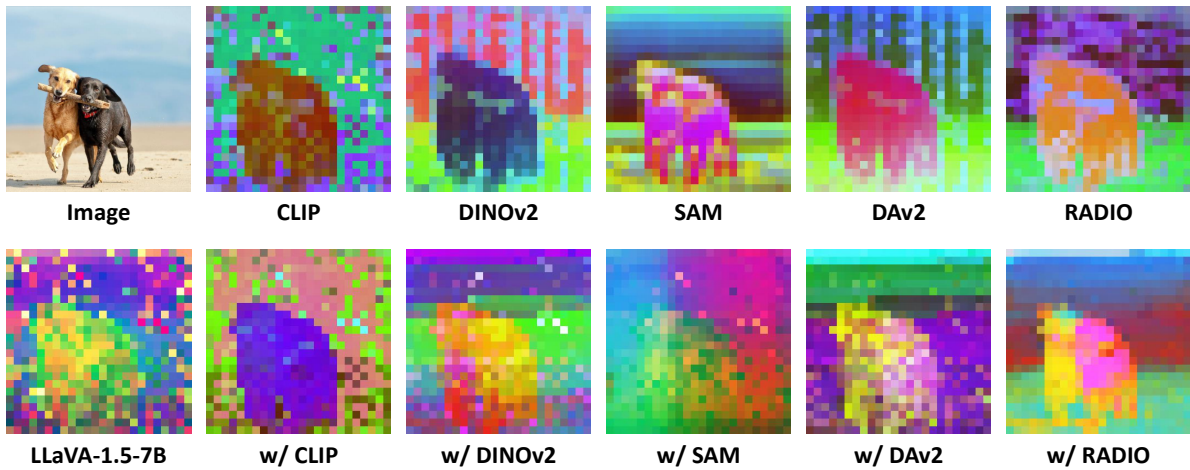


Figure 11: PCA visualizations of 16th layer visual representations aligned with different VFMs: CLIP (Radford et al., 2021), DINOv2 (Oquab et al., 2023), SAM (Kirillov et al., 2023), DAv2 (Yang et al., 2024), and RADIO (Heinrich et al., 2025).

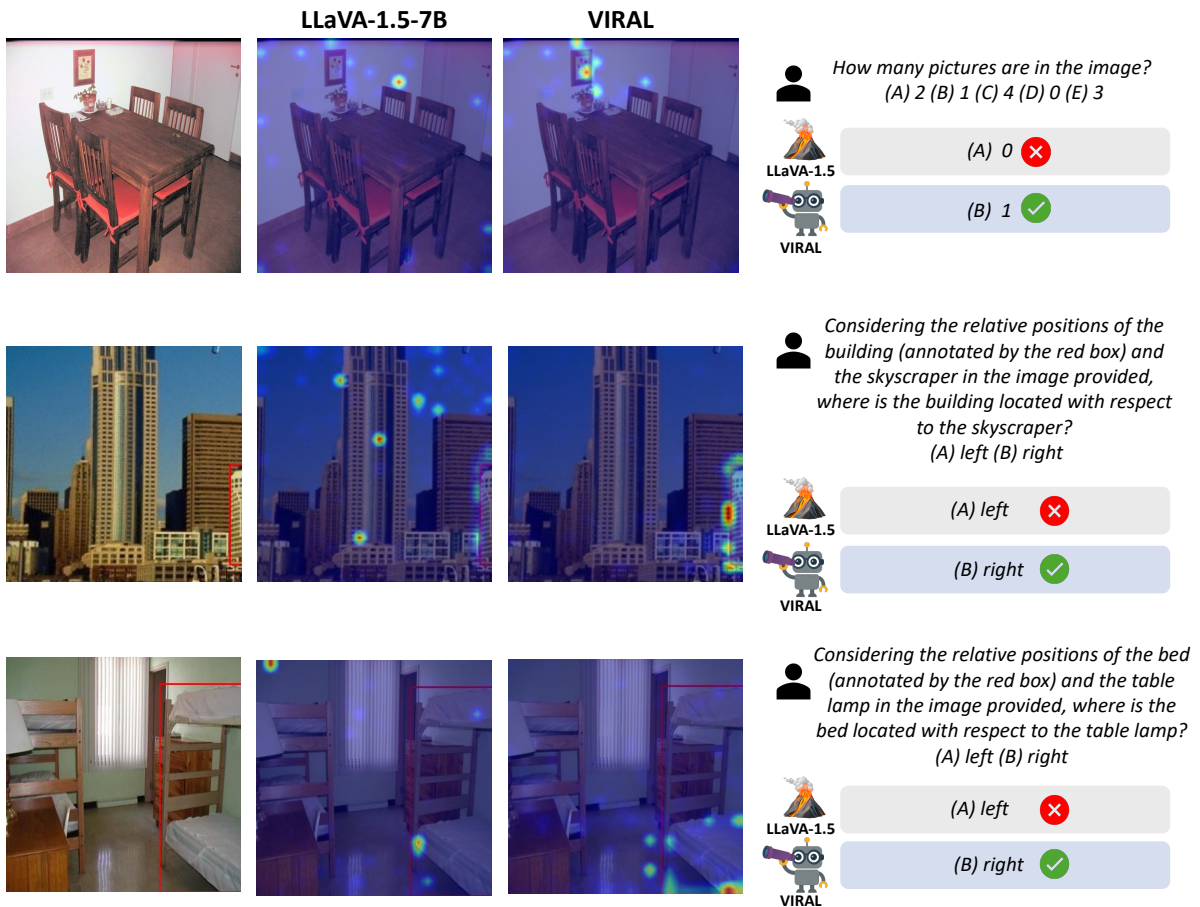


Figure 12: Cross-attention map comparison for vision-centric tasks.