Towards Mitigating Hallucinations in Large Vision-Language Models by Refining Textual Embeddings

Aakriti Agrawal ¹, Gouthaman KV ², Rohith Aralikatti ³, Gauri Jagatap ², Jiaxin Yuan ¹, Vijay Kamarshi ², Andrea Fanelli ², Furong Huang ^{1,4}
Correspondence: agrawal5@umd.edu ¹

¹University of Maryland, ²Dolby Laboratories, ³Hilabs, ⁴Capital One

Hallucinations in Large Vision-Language Models (LVLMs) remain a persistent challenge, often stemming from inadequate integration of visual information during multimodal reasoning. A key cause is the model's over-reliance on textual priors and underutilization of visual cues, leading to outputs that are linguistically fluent but visually inaccurate. For example, given an image of an empty kitchen countertop, an LVLM might hallucinate a "bowl of fruit" or "cup of coffee," relying on language associations rather than visual evidence. Most LVLMs incorporate visual features by appending them to the input stream of a pre-trained LLM and training on large-scale vision-language datasets. Our systematic analysis reveals that this strategy often leads to over-dependence on textual information due to the inherent bias of LLMs towards language-dominant representations. This imbalance skews attention towards the text over visual content, weakening the model's ability to ground outputs in visual inputs. To address this, we propose a simple yet effective visual feature incorporation method that encourages the model to learn visually-informed textual embeddings distinct from those of the base LLM and promotes a more balanced attention distribution. Experimental results across multiple hallucination benchmarks demonstrate that our method significantly reduces hallucinations and fosters more balanced multimodal reasoning. Notably, our approach achieves substantial gains, including +9.33% on MMVP-MLLM, +2.99% on POPE-AOKVQA, up to +3.4% on Merlin, and +3% on the hard-data split of HallusionBench.

1. Introduction

The advent of LLMs has transformed NLP, enabling tasks like machine translation, dialogue, and content generation with unprecedented accuracy and fluency. Building on this, Large Vision-Language Models (LVLMs) (Lin et al., 2023, Zhang et al., 2023a, Maaz et al., 2024) integrate visual and linguistic understanding in a unified framework, bridging text and visual modalities. This synergy has advanced tasks such as captioning (Chen et al., 2022), question-answering (Li et al., 2023a), multimodal retrieval (Lin et al., 2024), etc. As LVLMs advance, their adoption in domains such as healthcare, autonomous driving, and education is accelerating, expanding the role of AI in real-world multimodal applications.

Despite this progress, IVLMs remain prone to hallucinations—outputs that are fluent but not grounded in the visual input. These errors, which include fabricating or misinterpreting visual content, undermine reliability and hinder deployment in safety-critical settings. Fig. 1 illustrates failure cases from Video-LLaVA (Lin et al., 2023), a widely adopted LVLM. In one case, the model captions a scene as "moving it towards a person," despite the absence of both the person and action in the video—highlighting both object and action hallucination. Beyond these, LVLMs hallucinations span a wide range. Attribute hallucination involve assigning incorrect visual properties, such as describing "a red car" when the car is actually blue, or denying visible objects (e.g., "there are no chickens" despite their presence). Relation hallucination fabricate spatial or

≥bec		Medals wir	nter olympics 2	022			
at Lakes	Country	gold	Silver	Bronze	Total		
	Norway	9	4	2	15		A Service Serv
	Germany	12	10	5	27	11	
higan	China	16	8	13	37		*
	Attrib	ute Hal	lucinatio	on:		Object Hallucination:	Act
ew	Q: Acc	ording t	o the tab	le, is the		Q: what is a man walking down?	Q : v
ares	numbe	er of gol	d medals	of Norwa	ay 9	(GT) A: Road/trail	(GT
	in win	ter Olym	pics in 20	022?		(Video-LLava) A: A man is walking	(Vic

Figure 1: Hallucinations in Video-LLaVA (Lin et al., 2023).

contextual links, such as claiming "the person is jumping over the fence" when the person is merely standing beside it. In video-based settings, hallucinations may also include fabricated temporal dynamics, such as asserting "the person enters the room" when the individual never appears in the frame.

We hypothesize that a fundamental source of hallucinations in LVLMs arises from the prevailing architectural paradigm in which the visual information is appended as embeddings (typically extracted using frozen pre-trained visual encoders) to the input sequence of a pre-trained LLM (Fig. 2, top). This fused input is then pass to the model and fine-tune on large-scale vision-language datasets, such as image/video captioning, and VQA (Lin et al., 2023, He et al., 2024, Maaz et al., 2024), etc. While this approach offers modularity, data efficiency, and leverages the strong language generation capabilities of LLMs, it introduces a structural asymmetry: the LLM backbone, trained solely on text, remains inherently biased toward language-driven reasoning (An et al., 2025, Arif et al., 2025). As a result, during fine-tuning, the model may tend to fall back to text priors, under-utilizing the visual embeddings and treating them as secondary in the reasoning process. This modality imbalance may lead to a systematic misalignment between visual evidence and generated text, manifesting during inference as hallucinations: outputs that are linguistically coherent and semantically plausible, yet factually incorrect or unsupported by the visual input.

Motivated by this, we systematically investigate modality imbalance in IVLMs as a potential source of hallucinations, with a focus on the dominant practice of appending visual embeddings to the input textual tokens of pre-trained LLMs (Lin et al., 2023, He et al., 2024, Maaz et al., 2024). We use Video-LLaVA (Lin et al., 2023) as baseline due to its strong performance, modular design, and community adoption, making it well-suited for controlled analysis and evaluation. Our analysis reveals that the prevailing approach of simply appending visual embeddings to the textual input sequence causes the model to over-rely on language while under-utilizing visual information, thereby exacerbating hallucinations. This arises because the backbone LLM, optimized for text, disproportionately emphasizes textual tokens during self-attention operations within the transformer layers.

To address this modality imbalance, we propose a method that integrates visual information directly into text embeddings at the token level, enabling more balanced attention and cross-modal representations that better ground generation in visual input. By embedding visual semantics within the language representation, our method fosters balanced cross-modal reasoning, and reduces hallucinations. Extensive evaluation across multiple hallucination benchmarks demonstrates consistent and statistically significant gains, highlighting

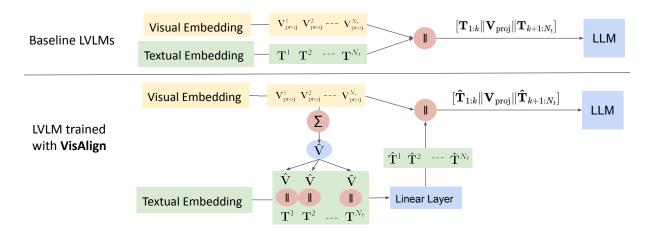


Figure 2: Top: Architecture of typical LVLMs like Video-LLaVA, which fuse language and vision embeddings by simple concatenation. **Bottom:** Our modified architecture with a **concatenation block that appends the averaged vision embedding to each token embedding, followed by a projection layer.** This encourages the model to learn visually informed textual embeddings and better attend to visual input during training.

the effectiveness and generalizability of our method.

2. Related Works

Large Vision-Language Models (LVLMs) extend pre-trained LLMs to handle visual inputs, typically by appending visual embeddings—extracted from frozen image or video encoders—to the language token sequence. This token-level fusion strategy enables architectural modularity and reusability of LLMs without major modifications. Notable models following this approach include LLaVA (Liu et al., 2024b), MiniGPT-4 (Zhu et al., 2023b), Video-LLaVA (Lin et al., 2023), Video-ChatGPT (Maaz et al., 2023), Bunny (He et al., 2024), and Video-LLaMA (Zhang et al., 2023a). Among these, Video-LLaVA has emerged as a foundational model due to its strong benchmark performance, open-source nature, and straightforward temporal extension via frame-wise token concatenation (Tang et al., 2025). Tang et al. (Tang et al., 2025) further identify Video-LLaVA as a key reference model that underpins many derivatives, including Video-ChatGPT (chat applications), Bunny (efficiency), and Video-LLaMA (fine-grained fusion).

Models like Flamingo (Alayrac et al., 2022) and BLIP-2 (Li et al., 2023a) use complex cross-attention to integrate modalities dynamically across transformer layers. Although more flexible, they incur higher computational costs and less modularity. Empirical results (Liu et al., 2024b) show simpler token-appending strategies often match or outperform these methods in accuracy and efficiency. For its simplicity, extensibility, and strong performance, we adopt Video-LLaVA (Lin et al., 2023) as our base model to investigate visual feature integration limitations, focusing on attention distribution, modality alignment, and hallucination.

Hallucination Detection and Mitigation in LVLMs is actively studied through a range of benchmarks designed to evaluate diverse hallucination types. POPE-AOKVQA (Li et al., 2023b) and NOPE (Lovenia et al., 2023) focus on object-level hallucinations, while MERLIN (Jing et al., 2023) examines factual consistency via atomic fact decomposition. MMVP-MLLM (Tong et al., 2024) and HallusionBench (Guan et al., 2024) probe model behavior under minimal semantic variation and cross-modal conflicts. Mementos (Wang et al., 2024a) targets temporal hallucinations in sequential visual reasoning. AMBER (Wang et al., 2023a) introduces a unified benchmark for evaluating both discriminative and generative hallucinations. Together, these datasets reveal a broad spectrum of hallucination phenomena—including object, action, attribute, relational, and

temporal inconsistencies—highlighting the complexity of achieving reliable visual grounding in LVLMs.

Several approaches have recently been proposed to mitigate hallucinations in LVLMs. M-HalDetect (Gunjal et al., 2023) introduces a dataset of hallucinated captions for training classifiers, while HaELM (Wang et al., 2023b) proposes a fine-tuning framework to distinguish hallucinated from faithful outputs. Reinforcement learning methods such as GAVIE (Liu et al., 2023) penalize ungrounded generations, and ALOHa (Petryk et al., 2023) leverages LLMs to detect hallucinated objects beyond fixed vocabularies. RLHF-based techniques (Sun et al., 2023) further enhance multimodal alignment. CLOCK (Biten et al., 2022) uses attention calibration during training. Inference-time strategies include visual-grounding-enhanced decoding via image descriptions (Ghosh et al., 2024), Instruction Contrastive Decoding (ICD) (Wang et al., 2024b), Self-Introspective Decoding (SID) (Huo et al., 2024), which verifies partial generations, and Visual Contrastive Decoding (VCD) (Leng et al., 2024), which re-ranks outputs to promote visual consistency. Together, these methods represent the current state of the art in hallucination mitigation.

Unlike prior approaches that rely on post-hoc corrections, inference-time heuristics, or hallucination-supervised fine-tuning, our method addresses hallucination proactively at the input representation level. By enriching textual embeddings with visual information, we encourage more balanced cross-modal attention and more effective utilization of visual cues—directly targeting the root cause of modality imbalance. This results in a more integrated, principled, and generalizable solution. Furthermore, our method not only performs well independently but also enhances existing techniques, as demonstrated in Appendix A, underscoring its broad applicability and complementary strengths.

3. Background

As noted above, we adopt the widely used Video-LLaVA as our baseline due to its pivotal role in advancing the field. Its canonical status and strong performance make it an ideal foundation for investigating modality imbalance, attention dynamics, and hallucination behaviors in LVLMs. This section formally outlines the architecture and training pipeline of Video-LLaVA (refer Figure 2 for an overview). It consists of the following components:

A frozen visual encoder to extract embeddings from the video (or image), the Video-LLaVA uses the pre-trained LanguageBind (Zhu et al., 2023a).

A projection layer that maps the visual embeddings into the textual (base LLM's) embedding space. The vision-language alignment is carried out via this projection layer. Formally, let $V \in \mathbb{R}^{N_v \times d_v}$ denote the visual embeddings, where N_v is the number of visual tokens and d_v is the visual embedding dimension. Output from the learnable projection layer $W_p \in \mathbb{R}^{d_v \times d_t}$ is denoted as:

$$V_{\text{proj}} = VW_p$$
, where $V_{\text{proj}} \in \mathbb{R}^{N_v \times d_t}$ (1)

where d_t is the LLM embedding dimension.

A backbone LLM: as mentioned above, LVLMs typically extend upon a pre-trained LLM. Video-LLaVA uses the pre-trained Vicuna-7b (Zheng et al., 2023).

The training involves the following two stages:

Pretraining: The visual encoder is frozen, and only the projection layer W_p will be trained. This stage aims to learn the visual-language alignment such that visual embeddings are interpretable by the base LLM.

Finetuning: In this stage, the entire pipeline including the LLM is trained end-to-end. The goal is to adapt the LLM to more effectively integrate and reason over visual embeddings alongside textual inputs, enabling

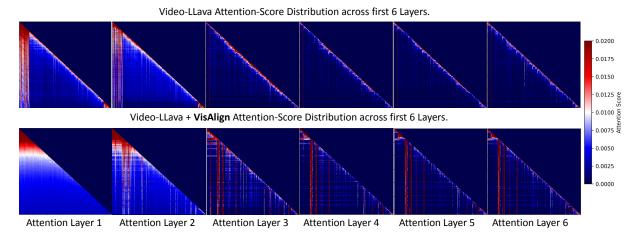


Figure 3: Attention score distributions across the first six attention layers of the baseline Video-LLaVA model (top row) and the VisAlign-enhanced model (bottom row). Video-LLaVA concatenates tokens in a fixed order: 35 initial text tokens, followed by 256 visual embeddings, and then the remaining text tokens. In each map, the x-axis denotes attended tokens (keys), and the y-axis denotes attending tokens (queries). Color intensity reflects attention weight: blue indicates low attention, red/white indicates high attention, and dark (near-black) regions indicate masked or negligible attention due to causal masking in autoregressive LVLMs.

visually-grounded language generation (e.g., captions or answers in VQA tasks).

4. Evaluating the Attention Score Distribution

Analyzing the attention score distribution across transformer layers provides insight into how information flows from lower to higher layers in LLMs. These scores reveal which tokens most influence the model's output and offer insight into its learning dynamics (Zhang et al., 2023b). Extending this analysis to LVLMs, we visualize the attention score distributions over both textual and visual tokens to better understand cross-modal interactions.

Figure 3 shows attention score distributions across multiple transformer layers in Video-LLaVA (Lin et al., 2023). In each heatmap, the horizontal axis represents Key tokens (tokens being attended to), and the vertical axis represents Query tokens (tokens performing attention). Color intensity encodes attention strength: cooler tones (e.g., blue) indicate lower scores, while warmer tones (e.g., red) and white indicate stronger attention. Nearly black regions show minimal or zero attention, often due to causal masking—a mechanism in auto-regressive LVLMs that prevents tokens from attending to future positions during decoding. This visualization qualitatively reveals how attention is distributed between visual and textual tokens across the network. Asymmetric or modality-skewed patterns highlight if the model overly favors one modality (typically text) at the expense of the other modality (visual), which can explain hallucination and grounding failures in multimodal tasks.

Figure 3 reveals a pronounced imbalance in how Video-LLaVA distributes attention between textual and visual tokens. In Layer 1 (top row, first plot), attention is heavily concentrated on the initial textual tokens (upper-left red region), sharply declines over the visual tokens, and rises again for the trailing textual tokens—a pattern consistent across layers. As defined in Eq. (1), the input sequence *X* follows a fixed order: initial textual tokens, followed by visual tokens, and then remaining textual tokens. This results in the model disproportionately focusing on textual tokens at both ends while under-attending to the visual tokens in between.

This asymmetric attention distribution reflects a modality bias rooted in the pre-trained base LLM, which was trained exclusively on text. During fine-tuning, the model relies heavily on linguistic priors and insufficiently leverages the visual embeddings provided by the frozen image or video encoder. This imbalance restricts the effective propagation and integration of visual signals across transformer layers, undermining robust visual grounding. Consequently, the model is prone to generate hallucinations—outputs that are fluent and semantically coherent but factually misaligned or unsupported by the visual input.

5. Improving Attention Score Distribution by Refining Textual Embeddings

We propose a simple yet principled approach, **VisAlign**, aimed at improving the attention score distribution across the visual and textual modalities. The underlying hypothesis is that encouraging a more balanced attention pattern, particularly by increasing attention to visual tokens, can help the model to better utilize visual information and reduce hallucinations caused by over-reliance on textual priors. VisAlign operates by refining textual embeddings through the integration of visual context prior to their input into the LLM. This alignment enables the model to extract and encode meaningful visual signals that might otherwise be underutilized due to the inherent linguistic bias of pre-trained LLMs. By fostering a more balanced and synergistic interaction between vision and language, VisAlign improves the utilization of visual information, without requiring architectural changes or external supervision.

As illustrated in Figure 2, VisAlign first applies average pooling on the projected visual embeddings $V_{\text{proj}} \in \mathbb{R}^{N_v \times d_t}$, resulting in the visual embedding vector $\hat{V} \in \mathbb{R}^{1 \times d_t}$:

$$\hat{V} = \frac{1}{N_v} \sum_{m=0}^{m=(N_v - 1)} V_{\text{proj}}[m]$$
 (2)

Next, we fuse \hat{V} with the text embeddings $T_{1:N_t} \in \mathbb{R}^{N_t \times d_t}$ via concatenation along the d_t dimension, yielding the fused embeddings T_V :

$$T_{V} = \left[T \mid \mid \hat{V} \otimes \mathbf{1}_{N_{t}} \right] \in \mathbb{R}^{N_{t} \times 2d_{t}}$$
(3)

Then, we apply a linear projection layer $W_d \in \mathbb{R}^{2d_t \times d_t}$ to map the fused representations T_V back to the original LLM embedding dimension d_t , producing the visually-grounded text token sequence, $\hat{T} = T_V W_d$ ($\in \mathbb{R}^{N_t \times d_t}$). Unlike the original textual tokens in base Video-LLaVA—which are derived solely from language embeddings—each token in \hat{T} encodes both visual and textual information. This enriched representation supports more effective cross-modal reasoning and visual grounding in downstream tasks.

Finally, we append \hat{T} to \hat{V} following the original concatenation strategy in Video-LLaVA (Eq. (1)):

$$\hat{X} = [\hat{T}_{1:k} || \hat{V} || \hat{T}_{k+1:N_t}]; \text{ where } \hat{X} \in \mathbb{R}^{(N_t + N_v) \times d_t}$$
(4)

The token sequence \hat{X} is then fed into the base LLM. It consists of visually grounded textual embeddings, followed by visual embeddings, and ends with the remaining grounded textual tokens.

Training Stages: We use the same datasets and training strategy as used in the baseline VideoLLaVA (Lin et al., 2023). In the *pretraining stage*, we train both the vision-language projection layer and the linear layer, while keeping the LLM frozen (refer to Figure 2 for an overview of the model). Whereas in the *finetuning stage*, we train the full model end-to-end, including the LLM, to allow for complete adaptation to the fused representations.

5.1. Attention score distribution with VisAlign

Figure 3 (bottom row) shows the attention distribution of Video-LLaVA trained with the VisAlign method. As illustrated, attention with VisAlign is more balanced and structured, spanning both visual and textual tokens throughout the sequence. Notably, the vertical attention bands are sharper and more frequent, indicating that the model consistently attends to specific visual regions or tokens that serve as semantic anchors across layers. Additionally, the smoother and more continuous diagonal gradients indicate that tokens attend not only to their local context but also capture long-range dependencies, reflecting a balanced and context-aware attention mechanism. In contrast, the top row (baseline Video-LLaVA) shows less coherent, more fragmented attention patterns. High attention is concentrated at the sequence boundaries, corresponding to textual token positions (Eq. (1)), revealing a strong bias toward language inputs. The lack of consistent vertical stripes further suggests limited focus on key visual elements, weakening the model's ability to maintain cross-modal grounding over time. Overall, attention in the baseline appears noisy and scattered across layers, indicating difficulty in forming stable associations between visual content and language queries.

These differences highlight VisAlign's effectiveness in improving the model's ability to integrate visual and textual modalities. By promoting more balanced attention, VisAlign improves focus on critical visual cues often overlooked by baseline Video-LLaVA, strengthening temporal and spatial coherence across the transformer layers and boosting overall visual information use.

6. Experiments and Results

We demonstrated that VisAlign improves attention distribution in Video-LLaVA. In this section, we evaluate its effectiveness in reducing hallucinations by comparing the original model with the VisAlign-augmented version across several benchmarks.

MMVP-MLLM (Tong et al., 2024) benchmark features carefully curated image pairs with highly similar CLIP embeddings, minimizing semantic divergence and emphasizing subtle visual distinctions. Each pair is accompanied by two binary-choice questions targeting fine-grained visual understanding. A model receives credit only if it answers both correctly, enforcing a strict criterion that rewards accurate visual grounding and penalizes reliance on language priors. This makes MMVP-MLLM particularly effective for evaluating hallucinations, as it compels models to rely on actual visual evidence rather than linguistic shortcuts or memorized associations.

The results in Table 1 show that Video-LLaVA enhanced with VisAlign achieves a substantial +9.33% improvement over the baseline. Since MMVP-MLLM is specifically designed to probe bias and hallucination in LVLMs by enforcing fine-grained visual discrimination under minimal semantic variance, this gain is especially significant. It demonstrates that VisAlign markedly strengthens the model's grounding in visual evidence rather than relying on linguistic priors, effectively reducing hallucinations and improving factual consistency. A qualitative comparison is presented in Figure 4. In the first example, the model must distinguish between two flame images—one round and the other elongated. The baseline Video-LLaVA incorrectly classifies both as "round," indicating over-reliance on memorized language patterns. In contrast, the VisAlign-enhanced model correctly differentiates the shapes, demonstrating stronger visual grounding. Similar improvements appear in other examples, underscoring VisAlign's effectiveness in reducing hallucinations and promoting accurate, cross-modal reasoning.

POPE (Li et al., 2023b) evaluates hallucinations through yes/no questions about object presence in images. "Yes" questions correspond to ground-truth objects, while "No" questions are adversarially crafted from the top-k most frequent object categories absent from the image. This setup exposes the model's reliance on language priors by testing its ability to reject visually unsupported but common objects. Following prior

	MMVP-MLLM	POPE A-OKVQA				
	Acc	Acc	P	R	F1	
Video-LLaVA	14	54.1	52.14	99.6	68.45	
+ VisAlign	23.33	57.09	53.9	98.33	69.63	

Table 1: Results on POPE A-OKVQA (Li et al., 2023b) & MMVP-MLLM (Tong et al., 2024). Acc: Accuracy, P:Precision, R:Recall, F1: F1 score.

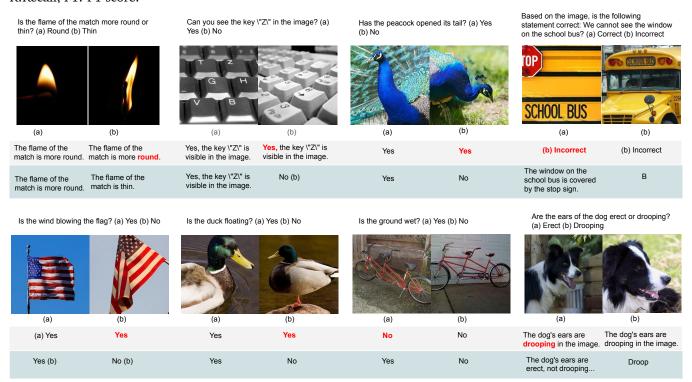


Figure 4: Qualitative results from the **MMVP-MLLM Benchmark**: Below each image, the baseline model's response is shown first, followed by the response from the model trained with *VisAlign*.

work (Villa et al., 2025), we focus on the most challenging setting: Adversarial SEEM from A-OKVQA, which applies SEEM-based object detection to A-OKVQA images. This subset probes whether models falsely affirm the presence of common yet incorrect objects, revealing object-level hallucinations driven by language bias. POPE thus offers a fine-grained, targeted measure of visual grounding, serving as a rigorous and complementary benchmark to evaluate VisAlign's effectiveness in reducing hallucinations.

Table 1 presents quantitative results on the POPE benchmark, where VisAlign consistently surpasses the baseline across key metrics, achieving a **2.99**% increase in accuracy, a **1.76**% boost in precision, and a **1.18**% gain in F1-score. The notable rise in precision indicates a significant reduction in false positives—hallucinated objects—while the improved F1-score reflects a more robust balance between precision and recall. These provide strong evidence that VisAlign effectively curtails predictions of frequent yet visually unsupported objects, thereby substantially enhancing object-level visual grounding. Supporting qualitative results in Fig. 5 further reinforce VisAlign's reliability in avoiding erroneous affirmations of absent objects, underscoring its critical role in advancing cross-modal integration and reducing hallucinations.

MERLIN Villa et al. (2023) evaluates factual consistency and visual grounding in LVLMs through fine-grained object existence verification. It employs a curated set of original and synthetically edited images to assess

	Curated Images								
	Pos-Orig	Pos-Edited	Neg-Orig	Neg-Edited					
Video-LLava	30.9	16.7	71.5	79.6					
VisAlign	34.3	20.3	72.7	83.0					
	Random Images								
	Pos-Orig	Pos-Edited	Neg-Orig	Neg-Edited					
Video-LLava	48.2	33.3	59.5	67.9					
VisAlign	48.6	36.7	60.1	71.3					

Table 2: Results (in %) on the Merlin benchamark Villa et al. (2023)). "Pos":Positive, "Neg":Negative.

		Ob	ject		Action						
Method	Acc.	Prec.	Rec.	F1	Acc.	Prec.	Rec.	F1			
Robotics Domain											
Video-LLaVA	8.27	16.55	12.40	13.46	5.53	6.99	11.30	8.40			
+VisAlign	9.40	19.20	13.61	15.16	6.50	9.60	10.76	9.45			
	Daily Life Domain										
Video-LLaVA	22.05	38.30	31.90	33.55	13.50	31.70	18.66	22.40			
+VisAlign	22.18	38.31	32.31	33.70	12.31	32.10	16.44	20.70			
Comics Domain											
Video-LLaVA	11.12	21.00	19.00	18.86	4.48	11.28	6.58	8.08			
+VisAlign	12.00	21.00	17.80	18.41	4.00	13.33	5.36	7.10			

Table 3: Results on Mementos (Wang et al., 2024a) across object and action hallucinations in three domains. whether models can accurately detect the presence or absence of objects. Our evaluation specifically targets a subset of MERLIN where an entire object category, limited to a single instance in the original image, has been removed in the edited version.

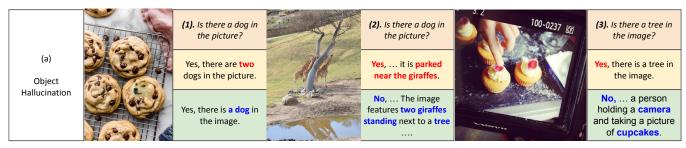
Table 2 presents quantitative results for both positive (object present) and negative (object removed) cases, evaluated under two distinct image sampling strategies. Across all configurations, VisAlign consistently outperforms the baseline, achieving significant improvements in accurately grounding object presence and absence. These results demonstrate VisAlign's superior capability to mitigate hallucinations by enhancing the model's sensitivity to subtle visual cues, thereby substantially improving visual fidelity and robustness in fine-grained, object-centric reasoning tasks.

Mementos (Wang et al., 2024a) evaluates sequential image reasoning in LVLMs across three domains: *Robotics*, *Comics*, and *Daily Life*. It rigorously test object and action hallucinations within dynamic visual contexts, emphasizing temporal coherence and object-behavior relationships. This makes Mementos especially valuable for assessing a multimodal model's ability to detect hallucinations while accurately understanding complex, evolving visual narratives.

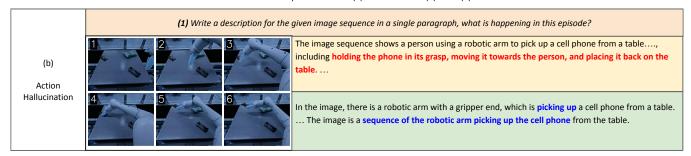
Table 3 shows significant improvements in the *Robotics* domain for both object hallucination (+1.13% accuracy) and action hallucination (+0.97% accuracy). These gains stem from the structured, goal-driven

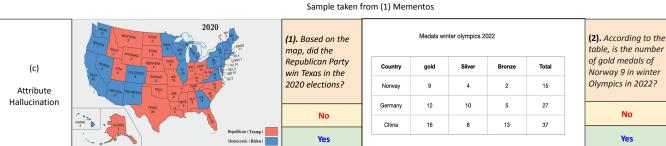
Visual Dependent						Visual Supplement				
Method	Figure	Ilusion	Math	OCR	Video	Chart	Map	OCR	Table	Average
Hard Da					ata Split					
Video-LLaVA	29.27	54.93	35.29	41.30	36.84	24.56	25.00	18.52	28.79	32.72
Video-LLaVA + VisAlign	34.15	49.30	37.25	45.65	36.84	21.05	28.12	33.33	34.85	35.61
	Easy Data Split									
Video-LLava	64.10	40.28	27.78	75.61	15.94	35.11	46.88	53.70	36.36	43.97
Video-LLava + VisAlign	53.85	36.11	37.04	53.66	36.23	25.95	48.44	50.00	28.57	41.1

Table 4: Category-wise results on the HallusionBench benchmark(Guan et al., 2024).



Samples taken from (1) HallusionBench (2) POPE (3) POPE







(d)

Relation Hallucination

Wisconsin also share a border with Lake Huron.

Wisconsin also share a border with Lake Huron.

Yes, according to the image, is New York the only U.S. state that shares a border with Lake Huron.

Wisconsin also share a border with Lake Huron.

Yes, according to the image, is New York the only U.S. state that shares a border with Lake Huron.

(a) Correct

(b) In the image, there is one daisy that is under the shadow of a taller daisy. (a) Correct (b) Incorrect

Yes, according to the image, New York is the only U.S state

Samples taken from (1) HallusionBench (2) MMVP

Figure 5: Qualitative examples from **POPE A-OKVQA**, **HallusionBench**, **MMVP**, and **Mementos** benchmarks illustrating various hallucination types. Input prompts are shown in orange, baseline Video-LLaVA outputs in yellow, and VisAlignenhanced outputs in green. VisAlign consistently improves performance across object, action, attribute, and relation hallucinations.

nature of robotic sequences, where predictable temporal patterns and clear visual cues enable VisAlign to maintain coherent attention over time and better align visual tokens with text, enhancing temporal reasoning of object states and behaviors. In contrast, improvements in the *Comics* and *Daily Life* domains are more modest, likely due to their greater visual and semantic complexity. Comics often use stylized, symbolic imagery and abstract narratives that disrupt typical visual-linguistic links, while Daily Life scenes involve high variability, subtle object transitions, and complex human actions that hinder consistent temporal alignment. In these unstructured contexts, VisAlign's attention calibration is limited by noisier, less reliable visual inputs.

HallusionBench (Guan et al., 2024) is a diagnostic benchmark assessing how parametric memory affects hallucinations in LVLMs. It categorizes questions into Visual-Dependent (VD), requiring visual input, and Visual-Supplement (VS), answerable using world knowledge or training data. VS questions evaluate the model's ability to resolve conflicts between visual input and parametric memory. The benchmark includes easy and hard splits, with the hard subset featuring human-edited images designed to create modality conflicts.

Table 4 shows an average improvement of about 3% on the challenging hard subset. Significant gains are seen in Visual-Dependent (VD) tasks, with improvements of 4.88%, 1.96%, and 4.35% in the "Figure," "Math," and "OCR" categories, respectively. Even larger gains occur in Visual-Supplement (VS) tasks, with 3.12%, 14.81%, and 6.06% improvements in "Map," "OCR," and "Table." These results are particularly notable because the hard subset contains human-edited images designed to conflict with common knowledge, forcing the model to rely on visual input rather than memorized facts. The gains indicate that VisAlign substantially improves the model's ability to ground predictions in visual evidence, reducing over-reliance on language priors. For example, Figure 5 (d)(1) shows a manipulated map where New York is falsely depicted bordering Lake Huron; while baseline Video-LLaVA hallucinates based on memorized geography, Video-LLaVA+VisAlign correctly interprets the altered visual context. Similarly, in (c)(2), a falsified medal count for Norway is accurately detected only by the VisAlign-enhanced model. These examples highlight VisAlign's effectiveness in enhancing visual grounding and mitigating hallucinations by improving sensitivity to subtle visual inconsistencies.

Summary: consistent improvements across all benchmarks demonstrate that refining attention score distributions effectively reduces hallucinations, enabling predictions grounded in visual evidence rather than memorized associations. Additional results in *Appendix A* show VisAlign's complementary gains when combined with state-of-the-art methods like VCD (Leng et al., 2024), as well as its performance on generic LVLM benchmarks and alternative baselines—underscoring its robustness and generalizability. Together, these findings reinforce VisAlign's effectiveness in enhancing visual grounding and reducing hallucinations in LVLMs.

7. Conclusion

We systematically analyze attention distributions in LVLMs concerning hallucinations—outputs lacking grounding in visual input. Our findings show that popular LVLMs like Video-LLaVA overemphasize text over visual information, increasing reliance on linguistic priors and hallucinations. To address this, we propose a simple yet effective method that enriches textual embeddings with visual cues, rebalancing attention and improving the model's use of visual information. This results in significantly reduced hallucinations and more semantically accurate, visually faithful outputs. We validate our approach across multiple challenging hallucination benchmarks, consistently achieving substantial improvements. We hope these insights inspire further research to better leverage visual data, reduce hallucinations, and enhance the reliability of multimodal reasoning in LVLMs.

8. Limitations and Future work

In this work, we identify an inherent bias in prevailing LVLM architectures toward the language modality, largely resulting from the common practice of simply appending visual embeddings to the input text sequence. To address this, we propose a simple yet effective method that refines textual embeddings by integrating average-pooled visual features. Our approach demonstrably improves visual grounding and significantly reduces hallucinations on established benchmarks. While average pooling offers a straightforward, robust, and efficient means of incorporating visual information, we believe that more sophisticated fusion methods could further enhance visual grounding and cross-modal alignment. Given that the primary focus of this work is to highlight the modality imbalance and its impact on hallucinations—and to show that refining textual embeddings with visual information mitigates this issue—we leave exploration of advanced fusion strategies for future work.

References

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- Wenbin An, Feng Tian, Sicong Leng, Jiahao Nie, Haonan Lin, QianYing Wang, Ping Chen, Xiaoqin Zhang, and Shijian Lu. Mitigating object hallucinations in large vision-language models with assembly of global and local attention. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 29915–29926, 2025.
- Kazi Hasan Ibn Arif, Sajib Acharjee Dip, Khizar Hussain, Lang Zhang, and Chris Thomas. Paint: Paying attention to informed tokens to mitigate hallucination in large vision-language model. *arXiv* preprint *arXiv*:2501.12206, 2025.
- Ali Furkan Biten, Lluis Gómez, Marçal Rusiñol, and Dimosthenis Karatzas. Let there be a clock on the beach: Reducing object hallucination in image captioning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2434–2443, 2022.
- Fu Chaoyou, Chen Peixian, Shen Yunhang, Qin Yulei, Zhang Mengdan, Lin Xu, Yang Jinrui, Zheng Xiawu, Li Ke, Sun Xing, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 3, 2023.
- Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual language-image model. arXiv preprint arXiv:2209.06794, 2022.
- Sreyan Ghosh, Chandra Kiran Reddy Evuru, Sonal Kumar, Utkarsh Tyagi, Oriol Nieto, Zeyu Jin, and Dinesh Manocha. Visual description grounding reduces hallucinations and boosts reasoning in lvlms. *arXiv* preprint *arXiv*:2405.15683, 2024.
- Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14375–14385, 2024.
- Prasanna Gunjal et al. M-haldetect: Detecting hallucinations in vision-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023.
- Muyang He, Yexin Liu, Boya Wu, Jianhao Yuan, Yueze Wang, Tiejun Huang, and Bo Zhao. Efficient multimodal learning from data-centric perspective. *arXiv preprint arXiv:2402.11530*, 2024.
- Jing Huo et al. Mitigating object hallucinations in large vision-language models via attention calibration. *arXiv* preprint arXiv:2502.01969, 2024.
- Long Jing, Zhe Wang, Yichen Zhang, Dacheng Tao, and Mingli Song. Faith: Faithful and informative textual hallucination detection in image captioning. In *Proceedings of the 2023 Conference on Computer Vision and Pattern Recognition*, pages 3456–3465. IEEE, 2023. doi: 10.1109/CVPR.2023.00345. URL https://openaccess.thecvf.com/content/CVPR2023/html/Jing_FAITH_Faithful_and_Informative_Textual_Hallucination_Detection_in_Image_Captioning_CVPR_2023_paper.html.

- Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. Mitigating object hallucinations in large vision-language models through visual contrastive decoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13872–13882, 2024.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023a.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023b.
- Bin Lin, Yang Ye, Bin Zhu, Jiaxi Cui, Munan Ning, Peng Jin, and Li Yuan. Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2023.
- Sheng-Chieh Lin, Chankyu Lee, Mohammad Shoeybi, Jimmy Lin, Bryan Catanzaro, and Wei Ping. Mmembed: Universal multimodal retrieval with multimodal llms, 2024. URL https://arxiv.org/abs/2411.02571.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 26296–26306, 2024a.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 26296–26306, June 2024b.
- Yuxiang Liu et al. Gavie: Grounded and verifiable image explanation. In EMNLP Findings, 2023.
- Holy Lovenia, Adji Bintang Wibowo, Krisna Kuntoro, Muhammad Firdaus, Radityo Eko Prasojo, Derry Tanti Suhendro, and Kurniawan Kurniawan. Nope: Evaluating and explaining negative object presence in image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1234–1243, 2023. doi: 10.1109/ICCV.2023.00123. URL https://openaccess.thecvf.com/content/ICCV2023/html/Lovenia_NOPE_Evaluating_and_Explaining_Negative_Object_Presence_in_Image_Captioning_ICCV_2023_paper.html.
- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023.
- Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL 2024)*, 2024.
- Nathan Petryk, Shikhar Sharma, Ali Furkan Biten, Lluis Gomez, Dimosthenis Karatzas, C V Jawahar, and Minesh Mathew. Aloha: Assessing language-only hallucinations in image captioning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6789–6798. Association for Computational Linguistics, 2023. doi: 10.18653/v1/2023.emnlp-main.567. URL https://aclanthology.org/2023.emnlp-main.567/.
- Hao Sun et al. Mmhal-bench: Multimodal hallucination benchmark for vision-language dialogue. *arXiv* preprint arXiv:2312.00704, 2023.

- Yunlong Tang, Jing Bi, Siting Xu, Luchuan Song, Susan Liang, Teng Wang, Daoan Zhang, Jie An, Jingyang Lin, Rongyi Zhu, et al. Video understanding with large language models: A survey. *IEEE Transactions on Circuits and Systems for Video Technology*, 2025.
- Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9568–9578, 2024.
- Andrés Villa, Juan Carlos León Alcázar, Alvaro Soto, and Bernard Ghanem. Behind the magic, merlim: Multi-modal evaluation benchmark for large image-language models. *arXiv preprint arXiv:2312.02219*, 2023.
- Andrés Villa, Juan León Alcázar, Motasem Alfarra, Vladimir Araujo, Alvaro Soto, and Bernard Ghanem. Eagle: Enhanced visual grounding minimizes hallucinations in instructional multimodal models. *arXiv* preprint arXiv:2501.02699, 2025.
- Xiyao Wang, Yuhang Zhou, Xiaoyu Liu, Hongjin Lu, Yuancheng Xu, Feihong He, Jaehong Yoon, Taixi Lu, Gedas Bertasius, Mohit Bansal, et al. Mementos: A comprehensive benchmark for multimodal large language model reasoning over image sequences. *arXiv preprint arXiv:2401.10529*, 2024a.
- Yixin Wang, Yuxiang Liu, Chunyuan Chen, Zhe Wang, Shuchang Yan, et al. Mitigating hallucinations in large vision-language models with instruction contrastive decoding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024b.
- Zhe Wang et al. Amber: A benchmark for evaluating hallucinations in multimodal models. *arXiv preprint arXiv:2310.12114*, 2023a.
- Zhe Wang et al. Haelm: Hallucination evaluation for large multimodal models. *arXiv preprint arXiv:2305.19162*, 2023b.
- Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023a.
- Qingru Zhang, Chandan Singh, Liyuan Liu, Xiaodong Liu, Bin Yu, Jianfeng Gao, and Tuo Zhao. Tell your model where to attend: Post-hoc attention steering for llms. *arXiv* preprint arXiv:2311.02262, 2023b.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623, 2023.
- Bin Zhu, Bin Lin, Munan Ning, Yang Yan, Jiaxi Cui, HongFa Wang, Yatian Pang, Wenhao Jiang, Junwu Zhang, Zongwei Li, et al. Languagebind: Extending video-language pretraining to n-modality by language-based semantic alignment. *arXiv* preprint *arXiv*:2310.01852, 2023a.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023b.

A. Appendix

Comparison with existing hallucination mitigation approaches: In the main paper, we showed that VisAlign significantly reduces hallucinations in Video-LLaVA by improving the attention score distribution across visual and textual modalities. In this section, we extend our analysis by comparing VisAlign with other state-of-the-art (SOTA) hallucination mitigation methods. As noted in the Related Work section (2), inference-time strategies currently represent the leading approaches for mitigating hallucinations. These methods intervene during the decoding stage to guide the model toward generating outputs that are more aligned with the visual input.

We focus on Visual Contrastive Decoding (VCD) (Leng et al., 2024), a strong inference-time SOTA method. VCD introduces a contrastive re-ranking mechanism, wherein multiple candidate responses are sampled from the model and scored based on both linguistic likelihood and visual alignment. This alignment is computed using a cross-modal similarity function that penalizes syntactically fluent yet visually inconsistent outputs. By re-ranking candidates, VCD encourages the model to favor generations that are both semantically coherent and grounded in the visual input—effectively reducing hallucinations without additional fine-tuning. While model-agnostic and lightweight, such inference-time methods complement VisAlign, which proactively mitigates hallucinations by refining representations during training.

Table 5 compares the effectiveness of VisAlign and VCD applied to Video-LLaVA, both individually and combined. VisAlign outperforms VCD alone, notably improving accuracy (54.5 57.09) and F1-score (68.6 69.63). While VCD delivers incremental gains by refining output selection during inference, VisAlign achieves more substantial improvements by addressing modality imbalance during training. When combined, the two methods yield the best overall performance, further boosting accuracy to 58.8 and F1-score to 70.04, demonstrating their complementary strengths.

These results underscore VisAlign's orthogonality to inference-time techniques like VCD, allowing it to enhance performance without interference. They also highlight its strong generalizability—VisAlign's benefits persist even when integrated with other hallucination mitigation strategies, showcasing its robustness across diverse settings.

Model	Acc	Precision	Recall	F1-Score
Video-LLaVA	54.1	52.14	99.6	68.45
+ VCD (Leng et al., 2024)	54.5	52.38	99.39	68.6
+ VisAlign	57.09	53.9	98.33	69.63
+ VisAlign + VCD	58.8	55.03	96.33	70.04

Table 5: Comparison of baseline Video-LLava with different combination of hallucination mitigation approaches on POPE-AOKVQA. We evaluate the individual and combined effects of Visual Contrastive Decoding (VCD), and VisAlign. The combination of both yields the best performance, with VisAlign contributing more significantly to hallucination reduction than VCD alone.

Model	Existence	Count	Position	Color	Posters	Celebrity	Scene	Landmark	Artwork	OCR
Video-LLaVA	170	121.66	88.33	135	103.74	101.47	163	161	107	87.5
+VCD (Leng et al., 2024)	170	105.00	76.66	125	100.00	100.88	155.75	154.5	99.25	77.5
+VisAlign	190	131.66	53.33	148.33	103.06	78.24	151	125	94	87.5

Table 6: Comparison of baseline Video-LLava with different combination of hallucination mitigation approaches on **MME**.

Effect of VisAlign on generic LVLM benchmarks: In the main paper, we comprehensively evaluated

VisAlign's effectiveness in reducing hallucinations across multiple benchmarks, consistently demonstrating significant and robust improvements. Although our primary focus is on hallucination tasks, to further investigate VisAlign's broader impact, we also assess how it influences the baseline model's performance on generic vision-language understanding benchmarks.

To this end, we evaluate on the MME benchmark (Chaoyou et al., 2023), a widely adopted diagnostic suite designed to probe the general capabilities of LVLMs. MME includes various subcategories covering fine-grained visual understanding and textual grounding tasks, such as Existence, Count, Position, Color, Posters, Celebrity, Scene, Landmark, Artwork, and OCR. These categories span a range of difficulty, from low-level visual perception to high-level semantic reasoning, offering a comprehensive lens into overall model competency.

Table 6 reports category-wise performance comparing the baseline Video-LLaVA, VisAlign and VCD augmented versions. VisAlign significantly improves upon the baseline in several key subcategories that are sensitive to visual grounding, such as Existence (170 190), Count (121.66 131.66), and Color (135 148.33). These improvements align with the primary objective of VisAlign—mitigating hallucinations by enhancing the model's attention to visual evidence—demonstrating its positive influence on tasks that demand precise object recognition and attribute understanding. Moreover, in categories such as OCR and Posters, VisAlign preserves the same level of performance as the baseline, indicating that it does not compromise tasks unrelated to hallucination-prone scenarios. However, some categories—such as Position, Celebrity, Scene, Landmark, and Artwork—show drop in performance. These tasks often require fine-grained spatial reasoning or prior world knowledge, which may be subtly impacted by VisAlign's architectural shift toward reinforcing visual embeddings over memorized linguistic patterns. This suggests that while VisAlign strengthens core visual grounding, it may introduce minor trade-offs in more specialized or context-dependent tasks.

Another observation from Table 6 is that state-of-the-art hallucination mitigation methods like **VCD cause** a universal performance drop or yield no improvements across all MME subcategories. In contrast, VisAlign demonstrates a more favorable trade-off: while it introduces minor performance reductions in certain high-level categories, it provides targeted improvements in core grounding tasks without degrading overall reliability. This contrast highlights VisAlign's orthogonality to inference-time methods and its potential to improve multimodal reasoning in a more integrated and generalizable manner.

In summary, while VisAlign is primarily designed to mitigate hallucinations, it also brings positive side effects on general VLM tasks that benefit from stronger visual grounding. By enriching textual embeddings with visual information, VisAlign promotes faithful grounding in visual inputs and reduces over-reliance on language priors. Unlike inference-time methods like VCD—which often reduce performance on generic benchmarks—VisAlign improves internal representations, preserving or enhancing accuracy in key subcategories like Color, Count, and Existence. However, this stronger grounding can slightly reduce performance in tasks relying on memorized knowledge or abstract reasoning (e.g., Landmark or Celebrity), due to reduced influence from language-driven biases. This trade-off is expected and could potentially be mitigated by training on larger-scale multimodal datasets—an exciting direction for future work. Overall, VisAlign offers a principled, generalizable, and training-efficient approach to hallucination reduction while preserving broader multimodal capabilities.

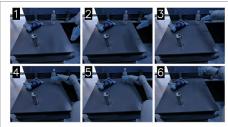
Performance on additional baselines: In the main paper, we demonstrated that VisAlign significantly reduces hallucinations in Video-LLaVA by improving attention distribution. To further validate the generality and robustness of VisAlign, we evaluate its effectiveness on another state-of-the-art LVLM, LLaVA 1.5 (Liu et al., 2024a). As shown in Table 7, VisAlign consistently enhances performance and reduces hallucinations when integrated into this baseline as well. These results highlight the broad applicability and effectiveness

Model	Acc	Precision	Recall	F1-Score
LLaVA1.5 (%)	69	62.23	97.66	76.02
+ VisAlign (%)	71	64	97.13	77.01

Table 7: Effects of VisAlign on the LLava1.5 baseline, on the POPE-AOKVQA benchmark

of the proposed approach across different LVLMs.

Additional qualitative results: Figure 6 (see next page) presents additional qualitative results on the Mementos dataset (Wang et al., 2024a). As illustrated, VisAlign enables the model to produce more visually grounded predictions and significantly reduces hallucinations compared to the baseline Video-LLaVA model.



Question: Write a description for the given image sequence in a single paragraph, what is happening in this episode?

Video-LLaVA: The image sequence shows a person using a vacuum cleaner to clean a table. In the first image, the person is holding the vacuum cleaner and preparing to clean the table. In the second image, the person is actively cleaning the table with the vacuum cleaner. In the third image, the person has finished cleaning the table and is holding the vacuum cleaner. In the fourth image, the person is emptying the vacuum cleaner into a trash can. In the fifth image, the person is putting the vacuum cleaner acts together. In the sixth image, the person is holding the vacuum cleaner again, possibly preparing to clean another area. Finally, in the seventh image, the person is emptying the vacuum cleaner into a trash can once again.

Video-LLaVA+VisAlign: In this image sequence, a robotic arm is seen picking up a bottle and placing it on a table. The scene is captured in four different frames, showing the progression of the action. The robotic arm is seen in different positions in each frame, indicating that it is in the process of picking up the bottle. The bottle is also appears in different positions in each frame, indicating that it is being picked up. The overall scene depicts a robotic arm picking up a bottle and placing it on a table.



Question: Write a description for the given image sequence in a single paragraph, what is happening in this episode?

Video-LLaVA: The image sequence shows a brown and white dog playing with a stuffed animal. In the first image, the dog is holding the stuffed animal in its mouth), while in the second image, the dog is biting the stuffed animal. The third image shows the dog holding the stuffed animal, and the fourth image shows the dog holding the stuffed animal in its mouth again. The fifth image shows the dog biting the stuffed animal, and the sixth image shows the dog chewing on the stuffed animal in its mouth. Overall, the sequence captures the playful and energetic behavior of the dog as it interacts with the stuffed animal.

Video-LLaVA+VisAlign: In this image sequence, there are two dogs, a brown and white dog and a black and white dog, playing together in a fenced area. The brown and white dog is biting the black and white dog's ear, and the black and white dog is trying to bite the brown and white dog is face. The two dogs are engaged in a playful interaction, with the brown and white dog being the aggressor and the black and white dog being the defender.



Question: Write a description for the given image sequence in a single paragraph, what is happening in this episode?

Video-LLaVA: In this image sequence, a group of people is gathered together, posing for a photo. They are standing in front of a building, and some of them are wearing backpacks. The group consists of at least 12 people, with some standing closer to the camera and others further back. The photo captures the essence of a group activity or event, with everyone smiling and enjoying the moment.

Video-LLaVA-VisAlign: In this image sequence, there is a man wearing a black jacket and a hat, standing in front of a group of people. The man is taking a selfie with a camera, capturing multiple pictures of himself and the group of people around him. The group consists of people of different ages and genders, and they are all posing for the camera. The man is taking the pictures from different angles and distances, capturing the whole group in each shot. The scene is a lively and fun atmosphere, with everyone enjoying themselves and having a good time.



Question: Write a description for the given image sequence in a single paragraph, what is happening in this episode?

Video-LLaVA: In this episode, a man is sitting on a chair in a bar, holding a bottle of beer. He is engaged in a conversation with another man who is standing nearby. The scene is depicted in three different drawings, each showing a different moment in the interaction between the two men. The first drawing shows the man sitting on the chair, while the second drawing shows the man standing next to the sitting man. The third drawing shows the man standing in front of the sitting man, indicating that the conversation is ongoing. The bar setting is evident from the presence of a bottle of beer and a chair, and the overall atmosphere appears to be casual and relaxed.

Video-LLaVA-VisAlign: In this comic strip, a man is seen stitting at a bar, holding a drink in his hand. Another man is standing behind the bar, pouring a drink into a glass. The man sitting is looking at the drink being poured, and appears to be enjoying the moment. The scene is set in a bar, with bottles and glasses visible in the background. The man pouring the drink is wearing a tie, and the man sitting is wearing a suit. The overall atmosphere of the scene is relaxed and casual, with a hint of humor as the man pouring the drink is exaggerating the pouring motion.

Figure 6: Qualitative results on the Mementos benchmark (Wang et al., 2024a). Text highlighted in red indicates hallucinated content, while text in blue shows the corresponding corrections.