ON ERRONEOUS AGREEMENTS OF CLIP IMAGE EM-BEDDINGS

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent research suggests that the failure of Vision-Language Models (VLMs) in visual reasoning could be attributed to the CLIP image encoder ambiguously encoding distinct images into embeddings with high cosine similarity, namely erroneous agreements. In this paper, we show that they are not the sole issue, as multimodal large language models (MLLMs) may extract distinct information even from image embeddings with high cosine similarities. On Subset A of the What'sUp benchmark, where the Left/Right image pairs are embedded by CLIP with average cosine similarity greater than 0.99, CLIP's performance is near random guess. In contrast, LLaVA-1.5-7B, which uses the same image encoder as CLIP, achieves nearly 100% accuracy. This discrepancy is also observed between LLaVA-1.5-7B and CLIP-like models on similar benchmarks. To investigate this performance gap, we conduct controlled experiments to test the effect of varying evaluation methods, training data, and language processing choices. We find that the CLIP image embeddings contain more extractable information than previously suggested, but it is likely obscured by the inadequate vision-language alignment of the CLIP's paradigm. Motivated by this observation, we reconsider the LLaVA-1.5 model on the MMVP benchmark, for which prior work showed that it could not distinguish image pairs with high cosine similarity. We observe a performance gain brought about by an alternative decoding algorithm, which attends more to visual input. Further, we show that the accuracy significantly increases if the model can take both images as input to emphasize their nuanced differences. Both findings indicate that LLaVA-1.5 did not utilize extracted visual information sufficiently. In conclusion, our findings suggest that while improving image encoders could benefit VLMs, there is room to enhance the models with a fixed image encoder through better strategies for extracting and utilizing visual information.

034

037

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

1 INTRODUCTION

Despite the rapid development and success of Vision-Language Models (VLMs), recent work pointed out that state-of-the-art VLMs (Radford et al., 2021; Zhai et al., 2023; Liu et al., 2024; Google, 2023a;b; OpenAI, 2023) still struggle with some simple visual reasoning tasks (Li et al., 2023d; Liu et al., 2023b; Tong et al., 2024c; Rahmanzadehgervi et al., 2024), where they were asked to answer questions about the images, such as recognizing shapes and describing the object relationships. These are basic tasks that VLMs should be able to solve before we deploy them to real-world scenarios like home robots responding to spoken or written commands.

Recent work argued that the pretrained CLIP image encoder (Radford et al., 2021), which serves as the "eyes" of many VLMs, could be the cause and cure for such visual shortcomings (Tong et al., 2024c). In these VLMs, any input image is first encoded by the CLIP image encoder and then used to calculate image-text similarity or as the input for a generative language model. Therefore, any deficiency of the CLIP image encoder propagates into the VLMs. The deficiency found by the authors is named *erroneous agreements*: Visually different images could be ambiguously encoded with high cosine similarity in the embedding space. They claimed this suggested information loss and caused VLMs' failure in relevant visual reasoning tasks, such as the MMVP benchmark (Tong et al., 2024c). This benchmark consists of selected, semantically distinct image pairs erroneously agreeing in the CLIP image embedding space, and CLIP-based VLMs failed to answer questions regarding the visual semantic difference better than random chance. This criterion is adopted in Taghipour et al. (2024) and similarly used in Tong et al. (2024b).

In this work, we provide evidence that VLMs face challenges beyond erroneous agreements. The 057 query-relevant visual information might still be present in the image embeddings despite the high cosine similarity, but a better strategy is required to pull it out. For instance, in the What'sUp benchmark (Kamath et al., 2023a) with paired, tightly controlled image pairs for evaluating VLM's 060 spatial reasoning ability, the average cosine similarity of image pairs on three out of four subsets 061 is greater than 0.95 in CLIP image embedding space, reaching the similarity threshold in Tong 062 et al. (2024c). While CLIP's accuracy in distinguishing these images is nearly random (about 50%), 063 LLaVA-1.5-7B (Liu et al., 2024) with the pretrained, frozen image encoder of CLIP-ViT-L/14-336px 064 still achieves beyond 80% in binary classification accuracies on all four subsets in What'sUp. Similarly, on the COCO-spatial and GQA-spatial benchmark used in Kamath et al. (2023a), LLaVA-065 1.5-7B surpasses CLIP-like models (including SigLIP (Zhai et al., 2023)) by a large margin. On the 066 more challenging MMVP and MMVP-VLM benchmark (Tong et al., 2024c), though its absolute 067 performance is poor, LLaVA-1.5 still outperforms CLIP-like models, showcasing a stronger ability 068 to extract information from given image embeddings. 069

What causes their discrepancy in extracting the given visual information? First, we unify the evaluation methods of CLIP and LLaVA and observe that the performance gap still exists. Then we decompose their difference into three parts: training data, language processing choice, and model paradigm (training and inference pipeline). Through ablation studies, we find that CLIP's failure is likely caused by the inadequate visual-language alignment of CLIP's paradigm. This also implies that the visual information extraction module in LLaVA-1.5, consisting of the two-layer MLP connector and the language model, adopts an inherently different mechanism from CLIP's paradigm.

077 The above results emphasize the importance of effective visual information extraction and highlight LLaVA-1.5's extracting ability. However, its poor performance on the MMVP benchmark remains a mystery. We look into its failure and provide insight into future directions in the discussion section. 079 To help LLaVA-1.5 keep the visual information during decoding, we try an alternative decoding algorithm, Multi-Modal Mutual-Information Decoding (M3ID) (Favero et al., 2024), leading to per-081 formance gain (+6%). We further find that visual nuances are often extracted and aligned with the correct semantics by LLaVA-1.5 rather than being discarded after visual encoding, but they did not 083 induce enough difference in outputs. To explore the amount of such visual formation, we reevaluate 084 LLaVA-1.5 with relaxed constraints, which allows for comparing the slight difference induced in the 085 outputs of two images. In this setting, its accuracy is significantly above random chance (+23.3%), while the result in the original one-image setting is just around random chance (+0.3%), suggest-087 ing insufficient visual information utilization in the original evaluation. In conclusion, despite the 088 erroneous agreements in the CLIP embedding space, visual nuances might still be extracted with improved strategies. This underscores the potential to enhance model performance by employing 089 better extraction and utilization techniques with the same pretrained image encoder. 090

- 091
- 092

2 Related Work

093 094

Benchmarking VLMs' visual reasoning ability. Vision reasoning tasks underline VLMs' visual 096 perception ability. Many recent challenging benchmarks on visual reasoning focus on assessing spe-097 cific abilities of current VLMs like compositionality (Winoground (Thrush et al., 2022), ARO (Yuk-098 sekgonul et al., 2023), SugarCrepe (Hsieh et al., 2024)), hallucination (POPE (Li et al., 2023d), HallusionBench (Liu et al., 2023a), and VHILT (Rawte et al., 2024)), distinguishing image pairs 100 (MMVP (Tong et al., 2024c)), spatial understanding of VLMs (What'sUp (Kamath et al., 2023a) 101 and Embspatial-bench (Du et al., 2024)), and core visual perception abilities (BLINK (Fu et al., 102 2024) for various aspects like visual correspondence and BlindTestbasic (Rahmanzadehgervi et al., 103 2024) for recognizing basic patterns). State-of-the-art VLMs often fail unexpectedly on simple test 104 cases, performing significantly worse than human accuracy or even random guess. For our discus-105 sion on erroneous agreements, we mainly consider the MMVP and What'sUp benchmark since their image pairs exhibit this property in the CLIP embedding space. Nevertheless, findings on these 106 benchmarks reveal the relationship between vision encoders and VLMs, supporting the broader goal 107 of enhancing VLMs for general visual reasoning.

124

125

126



Figure 1: An illustration of CLIP and LLaVA-1.5 model structures sharing the same pretrained image encoder with an example test case from the Left/Right subset of What'sUp benchmark. We find that the query-relevant nuances in the CLIP image embeddings may be extracted by LLaVA-1.5 despite erroneous agreements, and we note their performance gap on several similar benchmarks.

Exploring causes of the visual shortcomings of VLMs. Researchers are actively exploring the 127 root causes of VLMs' failures on benchmarks above, mainly from the model perspective. (1) Vision 128 modality. Tong et al. (2024c) argued that the flawed CLIP image encoder adopted by many VLM 129 architectures could lead to downstream failure because of erroneous agreements. (Chandhok et al., 130 2024) agreed that the image encoder is responsible for the information loss in spatial reasoning tasks 131 since CLIP's performance is quite low. However, we find evidence that erroneous agreements do not 132 necessarily lead to VLM's failure and that LLaVA is stronger at extracting information from simi-133 lar visual embeddings than CLIP. (2) Language understanding. In CLIP, Kamath et al. (2023b); 134 Tong et al. (2024b) found that its text encoder could also lose relevant information during encoding. 135 In multimodal LLMs (MLLMs), the language model might not timely terminate answer genera-136 tion (Yue et al., 2024), put false priority on the input text and format (Stan et al., 2024), or neglect 137 information like negation (Quantmeyer et al., 2024). Qiao et al. (2024) decoupled the perception stage and reasoning stage of VLMs and found that they are often limited by reasoning ability. In 138 this paper, we perform an ablation study on the text encoder and find that CLIP-like models still fail 139 when equipped with a stronger text encoder. (3) Vision-language alignment. Others discussed the 140 importance of modality alignment, such as visual grounding (Rajabi & Kosecka, 2023). From the 141 language model side, Ye et al. (2024) pointed out that the MLLMs might utilize multimodal spu-142 rious correlation in the training data due to the coarse-grained training objectives. Similarly, Yang 143 et al. (2024a) found that the model wrongly raised the probability of deceptive candidates. From the 144 image encoder side, Yang et al. (2024b) proposed the cross-modal Alignment and Correspondence 145 score of visual representations, which is linearly correlated to model performance. We abstractly 146 view the components of VLMs other than the vision part as visual information extraction and utilization module and demonstrate their different abilities. (4) Other factors. Apart from the model, 147 others looked into the problem with training data (Udandarao et al., 2024) or downstream tasks, such 148 as the hardness of the visual query (Zhang et al., 2024). 149

150 Improving the visual reasoning ability of VLMs. Following the observations about VLM's limi-151 tation, researchers mainly focused on improving the model structure with better image encoders and 152 vision-languageconnectors (Luo et al., 2024a; Zeng et al., 2021; Yao et al., 2024; Kar et al., 2024; Jiang et al., 2023; Zong et al., 2024; Xu et al., 2024; Tong et al., 2024a; Meng et al., 2024) or using 153 different training objectives with additional loss terms (Zhang et al., 2023; Zeng et al., 2024). From 154 the data-centric perspective, previous paper tried adding relevant instruction tuning data (Ranasinghe 155 et al., 2024; Chen et al., 2024), using long caption in pretraining (Zheng et al., 2024), hard negative 156 mining (Yuksekgonul et al., 2023; Paiss et al., 2023) or using synthesized images (Chatterjee et al., 157 2024; Jiao et al., 2024). 158

Post-training techniques are also explored to improve the performance of off-the-shelf VLMs. Some leveraged feedback from other models (Wang et al., 2024a; Luo et al., 2024b; Deng et al., 2024), various text or visual prompting methods (Wan et al., 2024; Lei et al., 2024; Wu et al., 2024), or multi-turn reflection (Huang et al., 2024; Wu et al., 2024; Wu & Xie, 2023; Kim et al., 2024b).



Figure 2: Example for evaluating CLIP-like models on What'sUp benchmark. For two-way evaluation, a test case consists of two similar images and two captions. The model chooses one caption for each image, and it gets one point in pair accuracy only if choosing correctly for both images.

Without external feedback or changing the task format, others developed probability-based output correction (Zhou et al., 2023), including hallucination-reducing decoding strategies (Chuang et al., 2023; Yang et al., 2024a; Kim et al., 2024a;c; Favero et al., 2024). Besides decoding, visual attention recalibration was proposed in response to the false priority put by VLMs (Woo et al., 2024). In this work, we achieve performance gain on MMVP through a decoding algorithm, M3ID (Favero et al., 2024), and test a new evaluation with relaxed constraints to show that current visual information utilization in LLaVA-1.5 is insufficient.

3 ERRONEOUS AGREEMENTS

We begin by introducing the task setup and the concept of erroneous agreements. Using a toy example, we demonstrate that information might still be extracted through alternative methods despite erroneous agreements. This is further validated through LLaVA-1.5's good performance on the What'sUp benchmark (Kamath et al., 2023a), showing that erroneous agreements are not the sole issue as it is possible for LLaVA-1.5 to extract distinct visual information from highly similar embeddings. We also notice a significant performance gap between LLaVA-1.5 and CLIP on What'sUp and across several other benchmarks.

3.1 TASK SETUP

200 This paper focuses on the setup in which VLMs are asked to choose from several captions based on a given image. For MLLMs, the image is accompanied by a question. Here, we use What'sUp 201 benchmark Kamath et al. (2023a), which was proposed for evaluating VLM's spatial reasoning 202 ability. Every test case includes four captions (e.g., "A dog left of/right of/on/under a table") and four 203 corresponding images photographed with minimal change except for the object spatial relationship. 204 For the convenience of calculating cosine similarity and comparing it to model performance, we 205 split each test case into two pairs: In the previous example, one pair consists of "A dog left of a 206 table" and "A dog right of a table" together with the ground truth images, and the other pair is the 207 remaining captions and images. This way, we get four subsets of the original benchmark. 208

For CLIP-like models, we calculate the matching score between images and texts. For CLIP with image encoder f_v and text encoder f_t , this is the cosine similarity between its image embeddings $f_v(\mathbf{v})$ and text embeddings $f_t(\mathbf{t})$, denoted as

212

178

179

180

181

182

183

185

186

187 188

189 190

199

- 213
- 214

$$S_C(f_v(\mathbf{v}), f_t(\mathbf{t})) = \frac{f_v(\mathbf{v})^\top f_t(\mathbf{t})}{||f_v(\mathbf{v})||||f_t(\mathbf{t})||}$$
(1)

As evaluation metrics, **pair accuracy** (Tong et al., 2024c; Kamath et al., 2023a) requires correct matching for both images, while the accuracy for two images independently is called **individual**

accuracy. An example from the What'sUp benchmark, together with the evaluation of CLIP, is
 shown in Figure 2.

The concept of erroneous agreements stems from this evaluation setup for CLIP. Specifically, if erroneous agreement happens for two different images v_1 and v_2 , then

$$_{C}(f_{v}(\mathbf{v_{1}}), f_{v}(\mathbf{v_{2}})) > \tau \tag{2}$$

where τ is a chosen threshold near 1 (e.g., $\tau = 0.95$ is used in Tong et al. (2024c)). Intuitively, when the cosine similarity is high, the two image embeddings point in nearly the same direction, and they will be close in Euclidean distance after l_2 -normalization. From the view of captions, they will result in highly similar image-text matching scores with any caption. From the image side, the margin is small, so adding a slight noise in either embedding could reverse the preference over the captions, which might result from a small perturbation in either v_1 or v_2 . Hence, the difference between v_1 and v_2 cannot be stably extracted by the CLIP model, and the relevant information seems lost. If this suggests the CLIP image encoder is "blind," this will also undermine VLMs that use it as "eyes."

This intuition is supported by the results on the MMVP benchmark designed to include image pairs with a cosine similarity greater than 0.95 for CLIP embeddings but less than 0.6 for DINOv2 embeddings, along with the MMVP-VLM benchmark for CLIP-like models (Tong et al., 2024c). LLaVA-1.5 and CLIP perform close to random chance on these two benchmarks, respectively. The authors also demonstrated a correlation between CLIP model accuracy and LLaVA-1.5's accuracy on different visual patterns.

236 237

241

221

3.2 DO ERRONEOUS AGREEMENTS MEAN BLINDNESS?

We note that cosine similarity does not depict all aspects of vector pairs. One criticism of it as the
 similarity metric is that it only captures the linear relationship of vectors. As an example, consider
 the following image embeddings

$$f_v(\mathbf{v_1}) = [10, 11, 12]^{\top}, f_v(\mathbf{v_2}) = [12, 11, 10]^{\top}$$

While $S_C(f_v(\mathbf{v_1}), f_v(\mathbf{v_2})) > 0.989$, Spearman's rank correlation coefficient can tell their sharp difference: $\rho = -1$, showing that their order information is fully opposed. Therefore, the difference in visual inputs might still be extracted through other means when erroneous agreements occur.

We show that this scenario happens in experiments. In many VLMs using CLIP image encoder as 246 their "eyes," the output score is nonlinear, different from CLIP-like models. For instance, in LLaVA-247 1.5 (Liu et al., 2024), the CLIP image embeddings first pass through a two-layer MLP and are then 248 used as input tokens for the transformer, Vicuna-1.5, which yields the token probability determining 249 the model response. We evaluate LLaVA-1.5-7B using the pretrained weights and design the ques-250 tion format. (An illustration and an example are in Figure 1, and more details are in Appendix A.1.) 251 The results are shown in Table 1. Despite the high cosine similarity, LLaVA-1.5-7B's individual 252 accuracy and pair accuracy are both quite high, showing that it can extract and align query-relevant 253 information from image embeddings and produce the correct answer. In other words, erroneous 254 agreements do not contribute to the failure of VLMs on their own.

255 Apart from this two-way evaluation, in Table 2, we report the results of the original evaluation, 256 which is a four-way classification for each image. Besides, we include the results on COCO-spatial 257 and GQA-spatial used in Kamath et al. (2023a) also for evaluating VLM's spatial reasoning ability. 258 These benchmarks are in the format of an image paired with two captions differing only by a prepo-259 sition. On all these benchmarks, LLaVA-1.5-7B wins CLIP-like models by a large margin, even 260 compared with the best model XVLM-COCO (Zeng et al., 2021) reported in the paper. We also find 261 that this performance gap relative to CLIP generalizes to some other MLLMs with different scales and language models in Appendix B.5. 262

To see if LLaVA-1.5-7B shows better extraction ability on tasks other than recognizing spatial relationships, we compare it and CLIP on MMVP and MMVP-VLM (Tong et al., 2024c). There was no direct comparison in the original paper: The MMVP benchmark is not in CLIP's format, while the MMVP-VLM benchmark is incompatible with MLLMs. So, we manually convert them into suitable formats without changing the content, and the evaluation of CLIP on MMVP-VLM is changed to the method described in Section 3.1 accordingly. The results are shown in Table 3. Although their absolute accuracy is low, there is a clear performance gap between CLIP-ViT-L/14-336px and LLaVA-1.5-7B with the same image encoder.

281

284

285

286 287

289

291

293

295

296 297

298

299

300

301

271 272

Table 1: The average cosine similarity of CLIP-ViT-L/14-336px embeddings and results of LLaVA-1.5-7B model on four subsets of What'sUp. The individual accuracy and pair accuracy are in percentage points. The average cosine similarity of the CLIP-ViT-L/14-336px image embeddings for 273 image pairs is calculated for each category. 274

	What'sUp Subset A				What'sUp Subset B			
	Left/I	Right	On/Under		Left/Right		Front/Behind	
	Indiv.	Pairs	Indiv.	Pairs	Indiv.	Pairs	Indiv.	Pairs
CLIP-ViT-L/14-336px	49.0	1.9	61.7	23.3	54.9	10.8	51.5	7.8
LLaVA-1.5-7B	99.0	98.1	80.1	60.2	100	100	98.5	97.1
Avg. Embedding Cosine Sim.	0.9	95	0.9	71	0.9	55	0.9	02

Table 2: Results of varied vision-language models on What'sUp, COCO-spatial, and GQA-spatial benchmark. We test the models on the four-way classification of each image. "Set of 4" is the correctness for all four images in a set.

	What	'alla G	hubcat A	What	'alle C	ubeat D	COCO	amoticl	COA	amotici
	wnat	sup s	Subset A	wnat	sup a	Subset B	COCO	-spatial	GQA-	spatial
	Indiv.	Pairs	Set of 4	Indiv.	Pairs	Set of 4	One-obj.	Two-obj.	One-obj.	Two-obj
CLIP-ViT-L/14-224px	26.7	1.0	0.0	25.7	1.5	0.0	49.1	50.2	46.0	48.1
CLIP-ViT-L/14-336px	28.9	1.0	0.0	27.2	1.0	0.0	48.9	51.1	46.6	49.1
SigLIP-ViT-L/16-384px	26.7	0.0	0.0	28.7	2.0	0.0	50.3	48.6	47.8	48.7
XVLM-COCO	41.8	17.0	1.9	42.2	15.7	2.9	68.4	73.6	69.1	67.0
LLaVA-1.5-7B	62.1	41.3	14.6	74.0	61.8	23.5	96.0	82.3	96.0	90.7
Random chance	25.0	6.3	0.4	25.0	6.3	0.4	50.0	50.0	50.0	50.0

4 INVESTIGATE THE PERFORMANCE GAP

The above performance gap might be caused by various factors: evaluation methods, training data, language processing choice, and model paradigm. Firstly, we used cosine-similarity-based evaluation for CLIP-like models and model-response-based evaluation for LLaVA-1.5. Apart from this, CLIP-like models and LLaVA-1.5 were trained on different data, adopted different language processing techniques, and were in different VLM paradigms.

302 In this section, we design and conduct ablation studies to determine whether these factors contribute 303 to the failure of CLIP-like models. The ablation of evaluation methods is conducted first to determine 304 whether there is really a performance gap. Then, we control the training data and text encoder of 305 CLIP-like models to see if they cause the failure (See the illustration in Figure 3). 306

307 4.1 UNIFIED EVALUATION 308

One might first question the different evaluations performed on the CLIP-like models and LLaVA-310 1.5. For the former, the evaluation is numeric-value based, while the latter is judged by its out-311 put response as a conversation agent, following the practice in previous work (Tong et al., 2024c). 312 Hence, we test LLaVA-1.5 again using standard Multiple-Choice (MC) evaluation, where we rank the perplexity of options ("A" and "B") calculated based on the probability of output tokens. This is 313 similar to the CLIP-like models' evaluation, where the image-text matching scores are ranked. 314

315 In Table 3, we observe that MC evaluation yields similar results on the MMVP benchmark and even 316 better results on MMVP-VLM. Thus, this does not count for their performance discrepancy. The 317 increase in results on MMVP-VLM is possibly related to the hallucination of MLLMs, where they 318 tend to follow their language prior during long answer generation and might gradually "forget" the 319 visual input.

- 320
- 321 4.2 TRAINING DATA 322
- The web-crawled image-caption corpora used for pretraining models like CLIP generally contain 323 very few high-quality, unambiguous image-caption pairs with prepositions, as pointed out in Kamath



Table 3: Results of CLIP and LLaVA-1.5 on the MMVP and MMVP-VLM benchmark. The last row is the LLaVA-1.5-7B accuracy in the Multiple-Choice setting. Accuracy with an asterisk is obtained on a converted version of the original benchmark.

Figure 3: Illustration for the CLIP paradigm and the ablation studies.

et al. (2023a). Equipped with the pretrained CLIP image encoder, LLaVA-1.5 was finetuned on 349 LCS-558K, a subset from LAION-CC-SBU with BLIP captions, and Instruction-following Data 350 Mixture (hereafter referred to as DataMix-665K) (Liu et al., 2024). These datasets were carefully 351 curated and thus of higher quality than web-crawled data. Besides, they are more relevant to the 352 spatial reasoning tasks: They have around 13K samples with phrases containing "left" or "right," 353 accounting for around 1% of all data. This ratio is higher than that in LAION-2B (English), where 354 various prepositions, including other directions like top and bottom, represent less than 0.22% of the 355 training data (Kamath et al., 2023a). Hence, we hypothesized that LLaVA-1.5's visual information 356 extraction ability benefits from these data.

357 To check the effect of training data, we use LLaVA-1.5's training data to fine-tune CLIP-like models. 358 We convert both datasets to the image-caption format (More details in Appendix B.2). By default, 359 we lock the image encoder during finetuning for strict ablation. The results are shown in Table 4. 360 Finetuning on LLaVA-1.5's training data slightly improves CLIP's performance, but it does not 361 help SigLIP. Still, their accuracy is around random chance. On SigLIP, we try unlocking the image 362 encoder during finetuning, but this does not increase model performance notably either (See results in Appendix B.3). In Appendix B.4, we also explore the effect of high-quality data on the gap 363 between LLaVA-1.5 and another VLM paradigm, XVLM (Zeng et al., 2021), and find that data do 364 not explain it solely. This result aligns with the failure in previous work to enhance CLIP models significantly by finetuning them on a much larger, preposition-focused subset of LAION (Kamath 366 et al., 2023a). 367

368 One might argue that the CLIP training objective differs from LLaVA-1.5, heavily relying on negative samples beyond data quality, so CLIP's failure on the new dataset might be due to the lack 369 of corresponding negatives. Next, we check this by observing whether negative samples help the 370 CLIP-like models learn better. For this experiment, we focus on the model's ability to distinguish 371 "left" and "right" and use the Left/Right subsets as the benchmarks. We construct hard negative 372 captions by switching the related phrases to their opposite, e.g., replacing "on the left" with "on the 373 right." The loss objective changes accordingly, following the NegCLIP method (Yuksekgonul et al., 374 2023). 375

The results are shown in Table 5. This strategy does not increase the model performance consistently on the Left/Right subset, which is observed in Kamath et al. (2023a) as well. Likewise, we try unlocking the image encoder of SigLIP in this setting, which does not make a big difference (See

324

328

330

331

332

333

334 335 336

337 338

339

341

342 343

345

Table 4: Results of CLIP and SigLIP on What'sUp, COCO-spatial, and GQA-spatial benchmark after finetuning on LLaVA-1.5's training data.

	What	'sUp S	Subset A	What	'sUp S	Subset B	COCO	-spatial	GQA-	spatial
	Indiv.	Pairs	Set of 4	Indiv.	Pairs	Set of 4	One-obj.	Two-obj.	One-obj.	Two-obj.
CLIP-ViT-L/14-336px	28.9	1.0	0.0	27.2	1.0	0.0	48.9	51.1	46.6	49.1
+ finetuning	31.1	9.7	0.0	30.6	7.4	0.0	54.2	55.5	51.0	52.6
SigLIP-ViT-L/16-384px	26.7	0.0	0.0	28.7	2.0	0.0	50.3	48.6	47.8	48.7
+ finetuning	27.2	1.9	0.0	24.8	2.0	0.0	49.3	50.5	47.0	54.6
Random chance	25.0	6.3	0.4	25.0	6.3	0.4	50.0	50.0	50.0	50.0

Table 5: Two-way evaluation results of CLIP and SigLIP focusing on the Left/Right subsets of What'sUp, COCO-spatial, and GQA-spatial benchmark with or without substituted text encoder, after finetuning on LLaVA-1.5's training data with or without hard negative captions. After finetuning, the accuracies are still around or below random chance.

	What'sU	p Subset A	What'sUp	Subset B	COCO	-spatial	GQA-	spatial
	Indiv.	Pairs	Indiv.	Pairs	One-obj.	Two-obj.	One-obj.	Two-obj.
CLIP-ViT-L/14-336px	49.0	1.9	54.9	10.8	51.6	48.4	52.1	50.4
+ finetuning	50.5	2.0	53.9	5.9	49.9	53.4	49.1	53.8
+ neg. cap.	50.5	1.0	50.5	1.0	48.5	55.6	49.4	50.4
+ llm2vec, finetuning	50.0	1.0	49.5	0.0	48.5	48.0	50.0	53.8
+ llm2vec, neg. cap.	49.5	2.9	50.5	6.9	48.8	46.2	48.1	51.9
SigLIP-ViT-L/16-384px	50.0	1.9	51.5	5.9	48.7	50.2	51.2	47.0
+ finetuning	49.0	1.0	51.0	3.9	50.8	53.1	49.7	55.3
+ neg. cap.	50.0	0.0	50.0	0.0	50.5	53.8	51.0	48.1
+ llm2vec, finetuning	50.5	2.9	51.0	3.9	50.1	54.8	49.1	51.1
+ llm2vec, neg. cap.	50.5	1.0	51.0	3.9	50.0	47.7	48.8	50.0
Random chance	50.0	25.0	50.0	25.0	50.0	50.0	50.0	50.0

results in Appendix B.3). All the above results indicate that data alone are not to blame for the failure of CLIP-like models.

4.3 LANGUAGE MODEL CHOICES

Previous research suggested that the CLIP text encoder is "blind" too (Tong et al., 2024b; Kamath et al., 2023b; Yuksekgonul et al., 2023)– It struggles with capturing changed word orders, negation, and spatial or numerical details. On the other hand, LLaVA-1.5-7B employs a pretrained large language model (LLM), Vicuna-1.5-7B, which is supposed to be better than the CLIP text encoder at language reasoning, and the commonsense knowledge it learned during language modeling should benefit VLMs.

Is pretrained LLM the secret to the success of LLaVA-1.5? To answer this question, we perform further experiments on finetuning CLIP-like models using both the LLaVA-1.5 training data and a stronger text encoder. Since Vicuna-1.5-7B is a decoder-only language model, we utilize the LLaMA-2-7B-chat-hf-mntp checkpoint provided in Llm2vec (BehnamGhader et al., 2024), where LLaMA-2-7B-chat model was converted to a text encoder and showed excellent performance in encoding texts. We replace the original CLIP text encoder with this pretrained model and add a two-layer MLP connector on top of it to align its output dimension with the CLIP image encoder. Since this text encoder is well-trained, we freeze it during finetuning. To make a fair comparison, we lock the image encoder again and use the same connector design as the one used in LLaVA-1.5-7B with only differing widths.

We train the models in two settings: plain data and data with hard negative captions. The hard
negative captions are constructed in the same way as in Section 4.2. The results are shown in Table 5. Surprisingly, a strong text encoder does not help either. Like the practice in Section 4.2, we
try unlocking the image encoder in this setting (See Appendix B.3). We only observe a significant
increase in the pair accuracy on What'sUp when a strong text encoder, hard negative captions, and
unlocked image encoder are all used. Still, in this case, the individual accuracy and the accura-

Table 6: Results of LLaVA-1.5-7B with M3ID ($\alpha = 0.6, \lambda = 0.15$) using original evaluation on 433 MMVP benchmark, along with the results of other methods. 434

435		Indiv. Acc.	Pairs Acc.
436			
437	LLaVA-1.5-7B	61.7	25.3
438	w/ RP (Jiao et al., 2024)	_	27.3
139	w/ M3ID (Favero et al., 2024)	64.3	31.3
439	w/ DIVA (Wang et al. 2024a)	_	31.3
440			
441	Random chance	50.0	25.0
442			

443 cies for the COCO-spatial and GQA-spatial do not improve. Based on these results, we argue that 444 higher quality and relevant training data or stronger language models do not solely contribute to the 445 performance gap.

446 By controlling other factors, we suggest that differences in VLM paradigms may largely explain the 447 performance gap. One hypothesis is that CLIP-like models use a dot product for image-text align-448 ment during training and inference. For a given image embedding, every text embedding is linearly 449 projected into a spectrum [-1, 1] regarding their matching degree. While this multimodal contrastive 450 paradigm achieves great success in tasks like zero-shot classification, the learned alignment might 451 not effectively capture all correspondences between image and text for various downstream tasks. This hypothesis aligns with our analysis in Section 3.1 that different visual information extraction 452 strategies matter. 453

454 455

456

432

5 DISCUSSION

457 Although we find that visual information extraction methods matter a lot, and LLaVA-1.5 has a 458 stronger extraction ability on highly similar embeddings, its poor performance on the MMVP bench-459 mark remains unexplained (Tong et al., 2024c). In this section, we reconsider its failure and provide 460 insight into future improvements based on two findings in MLLMs: They might not attend enough 461 to the visual input, and the visual information is often aligned correctly but probably did not induce 462 enough differences in the output token probability.

463 464

465

5.1 ALTERNATIVE DECODING FOR LLAVA

466 Inspired by the findings in Section 4.1 that MLLMs might "forget" the visual input gradually, one 467 possible improvement is to "remind" MLLMs of them, magnifying the effect of visual input on language models. Multi-Modal Mutual-Information Decoding (M3ID) was designed for this purpose 468 on MLLMs like LLaVA (Favero et al., 2024). For token in each decoding step t, M3ID computes the 469 output probability with the image and without any input image, denoted as l_c and l_u , respectively. 470 The latter corresponds to the language prior. Then a correction term $(l_c - l_u)$ is added to l_c with 471 weight $\frac{1-\exp(-\lambda t)}{\exp(-\lambda t)}$ if the model is not highly confident with the next token $(\max_k (l_c)_k < \log \alpha)$. 472 This correction prevents the VLM from omitting the visual input and relying on the language prior. 473

474 We test this decoding strategy on the MMVP benchmark in the standard setting. In Table 6, this 475 method achieves the most gain (+6%) relative to the baseline LLaVA-1.5-7b. We note that this 476 surpasses some methods that modified the vision part, such as Libra (30.0 with a decoupled and 477 more complex vision system) (Xu et al., 2024) and is on par with I-MoF (31.3 with interleaved 478 CLIP and DINO features) (Tong et al., 2024c). This result suggests that LLaVA-1.5 did not attend to the visual input enough and thus might miss the key information for answering the query. A 479 similar finding was described through the interpretability perspective in Stan et al. (2024). 480

481

482 5.2 EVALUATION WITH RELAXED CONSTRAINTS

483

We look into the results of the MC evaluation and find that the output token probability often dif-484 fers for two images (e.g., compared with image 2, image 1 slightly prefers caption 1 more). Still, 485 the evaluation omits them since we always pick the caption with a higher probability for each im-

487	Table 7: Results of CLIP-ViT-L/14-336px and LLaVA-1.5-7B using original pair evaluation and
488	new evaluation with relaxed constraints on MMVP benchmark.

489		Original Evaluation	w/ Relaxed Constraints
490 491	CLIP-ViT-L/14-336px LLaVA-1.5-7B	14.0	64.0 73.3
492 493	Random chance	25.0	50.0
494			

age. Such differences show that the visual nuances are often extracted and aligned with the correct 495 semantics by LLaVA-1.5 rather than being discarded after visual encoding. 496

497 How many visual nuances are preserved and extracted by LLaVA-1.5? We explore this question by 498 testing the LLaVA-1.5 on a new evaluation pipeline with relaxed constraints. To catch the slight 499 difference in model output, similar to the MC evaluation, we calculate the model perplexity of two possible options. MC only uses the letters "A" and "B" when computing perplexity, but we use the 500 full option for perplexity computation, e.g., "(a) Open" and "(b) Closed" in the questions provided by 501 the original benchmark. Denote the perplexity of two options (normalized by the number of tokens) 502 given two images to be $ppl_{i1c1}, ppl_{i1c2}, ppl_{i2c1}, ppl_{i2c2}$, respectively. We consider the model to be 503 correct for this test case if they satisfy 504

486

506

 $\frac{ppl_{i1c1}}{ppl_{i1c1} + ppl_{i1c2}} > \frac{ppl_{i2c1}}{ppl_{i2c1} + ppl_{i2c2}},$

507 In other words, the model considers (image 1, caption 1) with (image 2, caption 2) more possible 508 matches than (image 2, caption 1) with (image1, caption 2). This way, we "force" the model to 509 output differently for two images in a pair, and thus, the random chance is 50%. Through this 510 comparison, we amplify the semantics induced by visual nuances. We also apply this evaluation on CLIP, replacing perplexity with cosine similarity. In Table 7, the new performance is significantly 511 higher than random chance as the baseline (+23.3%), compared with the pair accuracy under the 512 original evaluation (+0.3%). This means more visual information can be extracted from the image 513 embedding and aligned with the correct semantics than the original results suggested. 514

515 The possible reason why they failed to be extracted in the original setting is that the language model 516 did not fully utilize the image-induced semantics. Consequently, they failed to affect the output probability enough to produce the correct answer. Influenced by language prior and spurious corre-517 lation with irrelevant text tokens (Ye et al., 2024), it will probably output common answers or even 518 hallucinations. Hence, reaching this "upper bound" in the original evaluation requires the VLM to 519 utilize extracted visual nuances properly and balance it with language prior during generation. 520

- 521 522
- CONCLUSION 6
- 523

524 Our study questions the use of erroneous agreements to reflect CLIP image encoders' information loss or blindness and finds that they are not the sole cause of VLM failures. We show that the 525 amount of extracted visual information largely depends on the extraction strategy, which varies 526 widely across VLMs. LLaVA-1.5, with a stronger extraction ability, outperforms CLIP-like models 527 on our benchmarks. Our controlled experiments suggest that the key factor in their performance 528 discrepancy might lie in their paradigms. We believe the information loss of the image encoder 529 should be defined when conditioning on the VLM paradigm and possibly the downstream task. 530

Our results suggest there is still room to enhance VLMs with a fixed, pretrained image encoder. 531 While balancing the visual grounding ability and image-text correspondence (e.g., combining pop-532 ular visual representation learning models in various styles) could reach the best trade-off on the 533 curve for heterogeneous benchmarks, developing advanced methods for VLMs to extract and utilize 534 given visual information might shift the curve upwards. 535

536 Limitation. We view the VLMs abstractly and do not look into fine-grained details on how the visual information is extracted and leads to the model's output. We leave its dissection for future research. For ablation studies, we do not train CLIP or SigLIP models from scratch or use larger 538 batch sizes due to the limitation in computing resources, so the conclusion on the effects of different factors is restricted.

540 REFERENCES 541

550

567

- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani 542 Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei 543 Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. Openflamingo: An open-544 source framework for training large autoregressive vision-language models. arXiv preprint arXiv:2308.01390, 2023. 546
- 547 Parishad BehnamGhader, Vaibhav Adlakha, Marius Mosbach, Dzmitry Bahdanau, Nicolas Chapa-548 dos, and Siva Reddy. Llm2vec: Large language models are secretly powerful text encoders. arXiv 549 preprint arXiv:2404.05961, 2024.
- Shivam Chandhok, Wan-Cyuan Fan, and Leonid Sigal. Response wide shut: Surprising observations 551 in basic vision language model capabilities. arXiv preprint arXiv:2408.06721, 2024. 552
- 553 Agneet Chatterjee, Yiran Luo, Tejas Gokhale, Yezhou Yang, and Chitta Baral. Revision: Rendering 554 tools enable spatial fidelity in vision-language models. arXiv preprint arXiv:2408.02231, 2024. 555
- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. 556 Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In Proceedings 557 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 14455–14465, 558 2024. 559
- Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian 561 Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual 562 language-image model. arXiv preprint arXiv:2209.06794, 2022. 563
- Yung-Sung Chuang, Yujia Xie, Hongyin Luo, Yoon Kim, James Glass, and Pengcheng He. Dola: 564 Decoding by contrasting layers improves factuality in large language models. arXiv preprint 565 arXiv:2309.03883, 2023. 566
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, 568 Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023. 570
- 571 Ailin Deng, Zhirui Chen, and Bryan Hooi. Seeing is believing: Mitigating hallucination in large vision-language models via clip-guided decoding. arXiv preprint arXiv:2402.15300, 2024. 572
- 573 Mengfei Du, Binhao Wu, Zejun Li, Xuanjing Huang, and Zhongyu Wei. Embspatial-bench: Bench-574 marking spatial understanding for embodied tasks with large vision-language models. arXiv 575 preprint arXiv:2406.05756, 2024. 576
- 577 Alessandro Favero, Luca Zancato, Matthew Trager, Siddharth Choudhary, Pramuditha Perera, Alessandro Achille, Ashwin Swaminathan, and Stefano Soatto. Multi-modal hallucination con-578 trol by visual information grounding. In Proceedings of the IEEE/CVF Conference on Computer 579 Vision and Pattern Recognition, pp. 14303–14312, 2024. 580
- 581 Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A 582 Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but 583 not perceive. arXiv preprint arXiv:2404.12390, 2024. 584
- 585 Google. Bard, 2023a.
- 586 Google. Gemini, 2023b. 587
- 588 Cheng-Yu Hsieh, Jieyu Zhang, Zixian Ma, Aniruddha Kembhavi, and Ranjay Krishna. Sugarcrepe: 589 Fixing hackable benchmarks for vision-language compositionality. Advances in neural informa-590 tion processing systems, 36, 2024. 591
- Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, 592 Lei Cui, Owais Khan Mohammed, Barun Patra, et al. Language is not all you need: Aligning perception with language models. Advances in Neural Information Processing Systems, 36, 2024.

- 594 Dongsheng Jiang, Yuchen Liu, Songlin Liu, Xiaopeng Zhang, Jin Li, Hongkai Xiong, and Qi Tian. 595 From clip to dino: Visual encoders shout in multi-modal large language models. 2023. 596 Qirui Jiao, Daoyuan Chen, Yilun Huang, Yaliang Li, and Ying Shen. Img-diff: Contrastive data 597 synthesis for multimodal large language models. arXiv preprint arXiv:2408.04594, 2024. 598 Amita Kamath, Jack Hessel, and Kai-Wei Chang. What's" up" with vision-language models? in-600 vestigating their struggle with spatial reasoning. arXiv preprint arXiv:2310.19785, 2023a. 601 Amita Kamath, Jack Hessel, and Kai-Wei Chang. Text encoders bottleneck compositionality in 602 contrastive vision-language models. arXiv preprint arXiv:2305.14897, 2023b. 603 604 Oğuzhan Fatih Kar, Alessio Tonioni, Petra Poklukar, Achin Kulshrestha, Amir Zamir, and Federico 605 Tombari. Brave: Broadening the visual encoding of vision-language models. arXiv preprint 606 arXiv:2404.07204, 2024. 607 608 Junho Kim, Hyunjun Kim, Yeonju Kim, and Yong Man Ro. Code: Contrasting self-generated description to combat hallucination in large multi-modal models. arXiv preprint arXiv:2406.01920, 609 2024a. 610 611 Junho Kim, Yeon Ju Kim, and Yong Man Ro. What if ...?: Counterfactual inception to mitigate 612 hallucination effects in large multimodal models. arXiv preprint arXiv:2403.13513, 2024b. 613 614 Sihyeon Kim, Boryeong Cho, Sangmin Bae, Sumyeong Ahn, and Se-Young Yun. Vacode: Visual augmented contrastive decoding. arXiv preprint arXiv:2408.05337, 2024c. 615 616 Xuanyu Lei, Zonghan Yang, Xinrui Chen, Peng Li, and Yang Liu. Scaffolding coordi-617 nates to promote vision-language coordination in large multi-modal models. arXiv preprint 618 arXiv:2402.12058, 2024. 619 620 Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Mimic-it: Multi-modal in-context instruction tuning. 2023a. 621 622 Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A 623 multi-modal model with in-context instruction tuning. arXiv preprint arXiv:2305.03726, 2023b. 624 625 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: bootstrapping language-image 626 pre-training with frozen image encoders and large language models. In ICML, 2023c. 627 Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating 628 object hallucination in large vision-language models. arXiv preprint arXiv:2305.10355, 2023d. 629 630 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr 631 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer 632 Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13, pp. 740-755. Springer, 2014. 633 634 Fuxiao Liu, Tianrui Guan, Zongxia Li, Lichang Chen, Yaser Yacoob, Dinesh Manocha, and Tianyi 635 Zhou. Hallusionbench: You see what you think? or you think what you see? an image-context 636 reasoning benchmark challenging for gpt-4v (ision), llava-1.5, and other multi-modality models. 637 arXiv preprint arXiv:2310.14566, 2023a. 638 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction 639 tuning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recogni-640 tion, pp. 26296–26306, 2024. 641 642 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, 643 Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around 644 player? arXiv preprint arXiv:2307.06281, 2023b. 645 Gen Luo, Yiyi Zhou, Yuxin Zhang, Xiawu Zheng, Xiaoshuai Sun, and Rongrong Ji. Feast your 646
- 647 eyes: Mixture-of-resolution adaptation for multimodal large language models. *arXiv preprint arXiv:2403.03003*, 2024a.

- Tiange Luo, Justin Johnson, and Honglak Lee. View selection for 3d captioning via diffusion ranking. *arXiv preprint arXiv:2404.07984*, 2024b.
- Lingchen Meng, Jianwei Yang, Rui Tian, Xiyang Dai, Zuxuan Wu, Jianfeng Gao, and Yu-Gang Jiang. Deepstack: Deeply stacking visual tokens is surprisingly simple and effective for lmms.
 arXiv preprint arXiv:2406.04334, 2024.
- 654 OpenAI. Gpt-4v(ision) system card, 2023.

681

689

690

- Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel.
 Teaching clip to count to ten. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3170–3180, 2023.
- Yuxuan Qiao, Haodong Duan, Xinyu Fang, Junming Yang, Lin Chen, Songyang Zhang, Jiaqi Wang, Dahua Lin, and Kai Chen. Prism: A framework for decoupling and assessing the capabilities of vlms. *arXiv preprint arXiv:2406.14544*, 2024.
- Vincent Quantmeyer, Pablo Mosteiro, and Albert Gatt. How and where does clip process negation?
 arXiv preprint arXiv:2407.10488, 2024.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Pooyan Rahmanzadehgervi, Logan Bolton, Mohammad Reza Taesiri, and Anh Totti Nguyen. Vision
 language models are blind. *arXiv preprint arXiv:2407.06581*, 2024.
- 672 Navid Rajabi and Jana Kosecka. Towards grounded visual spatial reasoning in multi-modal vision language models. *arXiv preprint arXiv:2308.09778*, 2023.
 674
- Kanchana Ranasinghe, Satya Narayan Shukla, Omid Poursaeed, Michael S Ryoo, and Tsung-Yu
 Lin. Learning to localize objects improves spatial reasoning in visual-llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12977–12987, 2024.
- Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fahad S. Khan. Llava++: Extending visual capabilities with llama-3 and phi-3, 2024. URL https://github.com/mbzuai-oryx/LLaVA-pp.
- Vipula Rawte, Anku Rani, Harshad Sharma, Neeraj Anand, Krishnav Rajbangshi, Amit Sheth, and
 Amitava Das. Visual hallucination: Definition, quantification, and prescriptive remediations.
 arXiv preprint arXiv:2403.17306, 2024.
- Gabriela Ben Melech Stan, Raanan Yehezkel Rohekar, Yaniv Gurwicz, Matthew Lyle Olson,
 Anahita Bhiwandiwalla, Estelle Aflalo, Chenfei Wu, Nan Duan, Shao-Yen Tseng, and Vasudev
 Lal. Lvlm-intrepret: An interpretability tool for large vision-language models. *arXiv preprint arXiv:2404.03118*, 2024.
 - Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023.
- Ashkan Taghipour, Morteza Ghahremani, Mohammed Bennamoun, Aref Miri Rekavandi, Zinuo Li,
 Hamid Laga, and Farid Boussaid. Faster image2video generation: A closer look at clip image
 embedding's impact on spatio-temporal cross-attentions. *arXiv preprint arXiv:2407.19205*, 2024.
- Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5238–5248, 2022.
- Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha
 Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. arXiv preprint arXiv:2406.16860, 2024a.

SUCHYDAIIY TOUY, LYTN JOHEN, AUG JACOD SICHHAIGH, MANS-DIOGUCHY TAHUEN OF HUUHHHOGA	systems
⁷⁰³ with language models. Advances in Neural Information Processing Systems, 36, 2024b.	sjotenis

- 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*, pp. 9568–9578, 2024c.
- Vishaal Udandarao, Ameya Prabhu, Adhiraj Ghosh, Yash Sharma, Philip HS Torr, Adel Bibi, Samuel Albanie, and Matthias Bethge. No" zero-shot" without exponential data: Pretraining concept frequency determines multimodal model performance. *arXiv preprint arXiv:2404.04125*, 2024.
- David Wan, Jaemin Cho, Elias Stengel-Eskin, and Mohit Bansal. Contrastive region guidance: Improving grounding in vision-language models without training. *arXiv preprint arXiv:2403.02325*, 2024.
- Wenxuan Wang, Quan Sun, Fan Zhang, Yepeng Tang, Jing Liu, and Xinlong Wang. Diffusion
 feedback helps clip see better. *arXiv preprint arXiv:2407.20171*, 2024a.
- Ziqi Wang, Hanlin Zhang, Xiner Li, Kuan-Hao Huang, Chi Han, Shuiwang Ji, Sham M Kakade, Hao Peng, and Heng Ji. Eliminating position bias of language models: A mechanistic approach. *arXiv preprint arXiv:2407.01100*, 2024b.
- Sangmin Woo, Donguk Kim, Jaehyuk Jang, Yubin Choi, and Changick Kim. Don't miss the forest for the trees: Attentional vision calibration for large vision language models. *arXiv preprint arXiv:2405.17820*, 2024.
- Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms.
 arXiv preprint arXiv:2312.14135, 2023.
- Yixuan Wu, Yizhou Wang, Shixiang Tang, Wenhao Wu, Tong He, Wanli Ouyang, Jian Wu, and
 Philip Torr. Dettoolchain: A new prompting paradigm to unleash detection ability of mllm. *arXiv* preprint arXiv:2403.12488, 2024.
- Yifan Xu, Xiaoshan Yang, Yaguang Song, and Changsheng Xu. Libra: Building decoupled vision
 system on large language models. *arXiv preprint arXiv:2405.10140*, 2024.
- Dingchen Yang, Bowen Cao, Guang Chen, and Changjun Jiang. Pensieve: Retrospect-then-compare mitigates visual hallucination. *arXiv preprint arXiv:2403.14401*, 2024a.
- Shijia Yang, Bohan Zhai, Quanzeng You, Jianbo Yuan, Hongxia Yang, and Chenfeng Xu. Law of vision representation in mllms. *arXiv preprint arXiv:2408.16357*, 2024b.
- Huanjin Yao, Wenhao Wu, Taojiannan Yang, YuXin Song, Mengxi Zhang, Haocheng Feng, Yifan Sun, Zhiheng Li, Wanli Ouyang, and Jingdong Wang. Dense connector for mllms. *arXiv preprint arXiv:2405.13800*, 2024.
- Wenqian Ye, Guangtao Zheng, Yunsheng Ma, Xu Cao, Bolin Lai, James M Rehg, and Aidong
 Zhang. Mm-spubench: Towards better understanding of spurious biases in multimodal llms. *arXiv preprint arXiv:2406.17126*, 2024.
- Zihao Yue, Liang Zhang, and Qin Jin. Less is more: Mitigating multimodal hallucination from an eos decision perspective. *arXiv preprint arXiv:2402.14545*, 2024.
- Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and
 why vision-language models behave like bags-of-words, and what to do about it? In *The Eleventh International Conference on Learning Representations*, 2023.
- Yan Zeng, Xinsong Zhang, and Hang Li. Multi-grained vision language pre-training: Aligning texts with visual concepts. *arXiv preprint arXiv:2111.08276*, 2021.
- Yunan Zeng, Yan Huang, Jinjin Zhang, Zequn Jie, Zhenhua Chai, and Liang Wang. Investigating compositional challenges in vision-language models for visual grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14141–14151, 2024.

756 757 758	Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 11975–11986, 2023.
760 761	Jiarui Zhang, Jinyi Hu, Mahyar Khayatkhoei, Filip Ilievski, and Maosong Sun. Exploring perceptual limitation of multimodal large language models. <i>arXiv preprint arXiv:2402.07384</i> , 2024.
762 763 764	Le Zhang, Rabiul Awal, and Aishwarya Agrawal. Contrasting intra-modal and ranking cross- modal hard negatives to enhance visio-linguistic fine-grained understanding. <i>arXiv preprint</i> <i>arXiv:2306.08832</i> , 2023.
765 766 767 768	Kecheng Zheng, Yifei Zhang, Wei Wu, Fan Lu, Shuailei Ma, Xin Jin, Wei Chen, and Yujun Shen. Dreamlip: Language-image pre-training with long captions. <i>arXiv preprint arXiv:2403.17007</i> , 2024.
769 770 771	Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and Huaxiu Yao. Analyzing and mitigating object hallucination in large vision-language models. <i>arXiv preprint arXiv:2310.00754</i> , 2023.
772 773 774 775 776	Zhuofan Zong, Bingqi Ma, Dazhong Shen, Guanglu Song, Hao Shao, Dongzhi Jiang, Hongsheng Li, and Yu Liu. Mova: Adapting mixture of vision experts to multimodal context. <i>arXiv preprint arXiv:2404.13046</i> , 2024.
776 777 778	
779 780	
781 782 783	
784 785	
786 787	
788 789 790	
791 792	
793 794	
795 796 797	
798 799	
800 801	
802 803	
804 805 806	
807 808	
809	

836 837

838 839

845 846

847

811	Table 8: Question form	nats for different subsets.
812	Subset	Question
813	Subset	Question
814	What'sUp Subset A&B, Left/Righ	Is the (object 1) to the left of or to
815		the right of the (object 2)? Answer
816		left or right.
817	What'sUp Subset A, On/Under	Is the (object 1) on or under the (ob-
818		ject 2)? Choose from the two op-
819	William C. Land D. Frank (Dalling)	tions.
820	what sup Subset B, Front/Benind	1 Is the (object 1) in front of or benind
821		hind
822	COCO/GOA-spatial One obj	Is the (object 1) on the
823	cocorogri spanai, one obj.	(left/right/top/bottom) or on
824		the (right/left/bottom/top)? Give a
825		short answer.
826	COCO-spatial, Two obj.	Is the (object 1) (to the left of/to
827		the right of/above/below) a (object
828		2) or (to the right of/to the left
829		of/below/above) a (object 2)? Give
830		a short answer.
000	GQA-spatial, Two obj.	Is the (object 1) to the
001		(left/right/front/behind) of a (object
002		2) or to the (right/left/behind/front)
000		of a (object 2)? Give a short
834		answer.
835		

А **BENCHMARKS AND EVALUATIONS**

840 We use the public pretrained weights of LLaVA-1.5-7B (https://huggingface.co/ llava-hf/llava-1.5-7b-hf) for evaluation and use greedy encoding by default to ensure 841 reproducibility. We use OpenAI's pretrained CLIP-ViT-L/14-336px model, SigLIP-ViT-L/16-384px 842 pretrained on the WebLI dataset (Chen et al., 2022) provided in the OpenCLIP repository, and offi-843 cial pretrained XVLM-16M weight for both evaluation and finetuning. 844

A.1 EVALUATION ON WHAT'SUP

848 The What'sUp benchmark (Kamath et al., 2023a) contains 820 images of pairs of household ob-849 jects, 408 in Subset A and 412 in Subset B. We corrected the mislabeled images in the GitHub Issues and reevaluated the pretrained VLMs. For CLIP, SigLIP, and XVLM's evaluation, we use 850 the official code provided by the What'sUp benchmark's authors in https://github.com/ 851 amitakamath/whatsup_vlms. 852

853 For LLaVA-1.5, the questions used for evaluation are listed in Table 8. Then the question is concate-854 nated with the fixed prompt template ("USER: <image> \land n(question) ASSISTANT:"). Consider-855 ing the position bias in LLMs (Wang et al., 2024b), we exchange the position of two prepositions in the question with 50% probability on COCO-spatial and GQA-spatial benchmarks for fair results. 856 On the What's Up benchmark, the orders are always the same for two images. Then, we evaluate the 857 outputs by keyword matching since we observe that the output is quite structured. 858

859 The reason why we use different commands after the main question (e.g., "Answer left or right", 860 "Choose from the two options", and "Give a short answer") is that we found the LLaVA-1.5 model 861 sensitive to such command. We tried "Answer on or under" for the On/Under subset in What'sUp Subset A, and the model accuracy is quite low. This is one of its limitations that deserves future 862 research. However, we aim to show that LLaVA-1.5 can extract such information, so we use the best 863 prompt to showcase its ability.

867 868



- 871
- 872 873
- 874
- 875
- 876 877

878

A.2 EVALUATION ON MMVP AND MMVP-VLM

The MMVP benchmark contains 150 pairs of similar images, and the MMVP-VLM benchmark has 135 pairs of similar images, divided into nine categories. There is an overlap between the image pairs in these two benchmarks. An example of an image pair and the corresponding prompt for LLaVA-1.5 in MMVP are shown in Figure 4. We corrected the mislabeled images in the GitHub Issues and reevaluated the pretrained VLMs. Since MMVP is incompatible with CLIP and so is MMVP-VLM with MLLMs, we convert their questions manually. We attach these new questions to the supplementary material for reference.

Figure 4: Example test case and prompt for LLaVA-1.5 in MMVP benchmark.

USER: <image>

open or closed? (a) Open

(b) Closed

Are the butterfly's wings closer to being

Answer with the option's letter from the

given choices directly. ASSISTANT:

In the standard setting, we evaluate the correctness of the model response by human evaluation. Although the accuracy given by the GPT-4 evaluation was close to that of a human evaluation on average, we noticed that it is unreliable since it gave several wrong judgments. So, we evaluate the correctness of the answers manually to avoid models getting higher accuracy by cheating GPT-4. In the Multiple-Choice setting, we calculate and rank the perplexity of "A</s>" and "B</s>" given by the model.

In Section 5.2, we calculate the perplexity of two options (The two options are "(a) Open" and "(b) Closed" for the example in Figure 4). We also add the EOS ("</s>") to the end of these options and normalize the perplexity by their number of tokens.

- 895
- 896 897

899 900

901

902

903

B SUPPLEMENTARY EXPERIMENTAL DETAILS AND RESULTS

B.1 HYPERPARAMETERS

Our code is based on https://github.com/mlfoundations/open_clip. We finetune CLIP and SigLIP models for 5 epochs with a learning rate of 5e-6 on the combination of converted LCS-558K plus converted DataMix-665K. We use 50 steps of warmup and AdamW optimizer with a cosine-annealing learning rate schedule. The batch size is 512, and we train the models on 4 gpus.

904 905 906

907

B.2 LLAVA-1.5'S TRAINING DATA

We check the frequency of appearance of the following keywords in DataMix-665K and LCS-558K:
"on the left," "on the right," "to the left," "at the left," "at the right." In DataMix-665K,
there are 12957 instances with at least one of the key phrases, among which 12658 have a paired
image. For captions (ground truth answers), this number is 13473 since an instance is paired with a
multi-turn conversation. In LCS-558K, there are 560 such instances and captions since each instance
has only one question and one answer.

In our experiments in Section 4.2, LCS-558K was converted from image-text pair format to conversation format, so we revert this process by using ground truth answer as the caption. Since
DataMix-665K is in a multi-turn conversation format, we randomly pick one answer as the caption in each epoch. In Section 4.3, the new text encoder can encode long paragraphs, so we use the concatenation of all answers in the multi-turn conversation as the ground truth caption.

Table 9: Results of SigLIP-ViT-L/16-384px on the Left/Right subsets of What'sUp, COCO-spatial,
 and GQA-spatial benchmark after finetuning with image encoder unlocked on LLaVA-1.5's training
 data (LCS-558K + DataMix-665K) with constructed hard negative captions.

	What'sUp	What'sUp Subset A		What'sUp Subset B		-spatial	GQA-spatial	
	Indiv.	Pairs	Indiv.	Pairs	One-obj.	Two-obj.	One-obj.	Two-obj.
SigLIP-ViT-L/16-384px	50.0	1.9	51.5	5.9	48.7	50.2	51.2	47.0
+ finetuning	50.5	2.9	51.5	5.9	48.7	57.7	50.5	48.1
+ neg. cap.	50.0	3.9	47.1	2.0	52.3	47.0	51.8	52.7
Random chance	50.0	25.0	50.0	25.0	50.0	50.0	50.0	50.0

Table 10: Results of SigLIP-ViT-L/16-384px on the Left/Right subsets of What'sUp, COCO-spatial, and GQA-spatial benchmark. We substituted the text encoder to be Llama-2-7b-chat-hf-mntp, then finetuned the model with image encoder unlocked on LLaVA-1.5's training data (LCS-558K + DataMix-665K) with or without constructed hard negative captions.

	What'sUp	Subset A	What'sUp Subset B		COCO	-spatial	GQA-spatial	
	Indiv.	Pairs	Indiv.	Pairs	One-obj.	Two-obj.	One-obj.	Two-obj.
SigLIP-ViT-L/16-384px	50.0	1.9	51.5	5.9	48.7	50.2	51.2	47.0
+ finetuning	50.0	1.0	49.0	7.8	50.9	50.2	48.7	50.8
+ neg. cap.	56.3	26.2	55.4	25.5	50.8	48.4	46.7	53.4
Random chance	50.0	25.0	50.0	25.0	50.0	50.0	50.0	50.0

B.3 RESULTS OF UNLOCKING IMAGE ENCODER

We try unlocking the image encoder during finetuning on the SigLIP model. The results after finetuning with CLIP text encoder are in Table 9, and results with LLaMA-2-7B-chat-hf-mntp are in
Table 10. Interestingly, we observe a significant increase in pair accuracy on the What'sUp benchmark only when using hard negative captions and a strong text encoder while unlocking the image
encoder. Still, the individual accuracy remains low.

950 951

944

918

931

932

933

B.4 RESULTS OF FINETUNING XVLM

Observing the similar failure of the data-informed attempt, previous work concluded that even with relevant, high-quality data and hard negatives, denser supervision is likely required to let the model learn the basic spatial relations (Kamath et al., 2023a), as in XVLM (Zeng et al., 2021), a VLM with supervision at the bounding-box level. However, LLaVA does not incorporate downstream task-related inductive bias or denser supervision to achieve high accuracy, yet it beats XVLM finetuned on COCO (Lin et al., 2014) on the What'sUp benchmark.

We explore finetuning XVLM on LLaVA's training data based on their official code (https: //github.com/zengyan-97/X-VLM), but no improvement is observed in the results (the last two model rows in Table 11). The image encoder is locked during finetuning. We use both contrastive learning loss and image-text matching loss. The evaluation is performed through the image-text matching score. We finetune the XVLM-16M model for 5 epochs with a learning rate of 1e - 5 and a weight decay rate of 0.01. We use 10% steps of warmup and AdamW optimizer with a lambda learning rate schedule. The batch size is 128, and we train the model on 4 gpus.

965

967

966 B.5 RESULTS OF DIFFERENT MLLMS

Do our findings on LLaVA-1.5 in Section 3.2 generalize to other MLLMs? We evaluate four other
MLLMs using their officially released weights. First, we consider LLaMA-3-V-8B and Phi-3-V3.8B which have LLaVA-like architecture and use frozen CLIP-ViT-L/14-336px as the image encoder (Rasheed et al., 2024). For MLLMs with different architectures and training data, we use Otter-Image-MPT7B (Li et al., 2023b;a) with frozen CLIP-ViT-L/14 as the image encoder and

Table 11: Results of XVLM-16M on the Left/Right subsets of What'sUp, COCO-spatial, and GQAspatial benchmark on LLaVA-1.5's training data.

	What'sUp Subset A			What	'sUp S	ubset B	COCO	-spatial	GQA-spatial		
	Indiv.	Pairs	Set of 4	Indiv.	Pairs	Set of 4	One-obj.	Two-obj.	One-obj.	Two-obj.	
XVLM-16M	50.0	30.6	1.0	32.8	10.8	0.0	65.4	64.6	63.2	53.3	
+ finetuning	46.4	28.4	1.0	34.6	8.3	1.0	66.8	65.2	61.3	51.2	
Random chance	25.0	6.3	0.4	25.0	6.3	0.4	50.0	50.0	50.0	50.0	

Table 12: Results of CLIP-ViT-L/14-336px and MLLMs on four subsets in What'sUp. The individual accuracy and pair accuracy are in percentage points. The average cosine similarity of the CLIP-ViT-L/14-336px image embeddings for image pairs is calculated for each category.

	What'sUp Subse			et A Inder	W] Left/	nat'sU Right	Jp Subset B Front/Behind	
	Indiv.	Pairs	Indiv.	Pairs	Indiv.	Pairs	Indiv.	Pairs
CLIP-ViT-L/14-336px	49.0	1.9	61.7	23.3	54.9	10.8	51.5	7.8
LLaVA-1.5-7B	99.0	98.1	80.1	60.2	100	100	98.5	97.1
LLaMA-3-V-8B	90.3	80.6	57.8	20.4	71.1	46.1	69.1	41.2
Phi-3-V-3.8B	100	100	85.4	70.9	100	100	56.9	13.7
InstructBLIP-Vicuna-7B	50.0	1.9	93.7	87.4	50.0	0.0	50.0	5.9
Otter-Image-MPT7B	50.0	1.0	56.8	13.6	50.0	0.0	51.5	11.8
Avg. Embedding Cosine Sim.	0.9	95	0.9	71	0.9	55	0.9	902

InstructBLIP-Vicuna-7B (Dai et al., 2023) with frozen EVA-CLIP-ViT-G/14 (Sun et al., 2023). Otter adopts the OpenFlamingo (Awadalla et al., 2023) paradigm with a Perceiver resampler module on top of the frozen image encoder, and then sends the output of this module to the cross-attention layers of the language model. InstructBLIP employs the pretrained BLIP-2 (Li et al., 2023c) model, with a Q-Former and a fully connected layer as the vision-language connector between the frozen image encoder and the language model. Inside the Q-Former, image embeddings are used in cross-attention layers.

During evaluation, we find that all of these MLLMs are sensitive to the wording in the command part, so we try several commands and report the best results as we did for LLaVA-1.5. We keep all other settings the same as in Section A.1.

The results are shown in Table 12 and Table 13. For comparison, we also include the results of CLIP-ViT-L/14-336px and LLaVA-1.5. The good performance of LLaMA-3-V-8B and Phi-3-V-3.8B verifies that they can also extract distinct information from highly similar embeddings, though they are weak on some prepositions (Front/Behind for Phi-3-V-3.8B, and On/Under for LLaMA-3-V-8B). These results show that our findings generalize to these two MLLMs with different scales and language models.

On the other hand, with different architectures and training data, Otter and InstructBLIP still struggle
on this benchmark (except On/Under for InstructBLIP). Hence, MLLMs do not guarantee effective
extraction from frozen image encoder. Good design of MLLM architecture and curated training data
synergize to provide strong visual information extraction ability.

Table 13: Results of CLIP-ViT-L/14-336px and MLLMs on What'sUp (four-way classification) benchmark, COCO-spatial, and GQA-spatial. "Set of 4" is the correctness for all four images in a set.

	What'sUp Subset A			What'sUp Subset B			COCO -spatial		GQA-spatial	
	Indiv.	Pairs	Set of 4	Indiv.	Pairs	Set of 4	One-obj.	Two-obj.	One-obj.	Two-obj.
CLIP-ViT-L/14-336px	28.9	1.0	0.0	27.2	1.0	0.0	48.9	51.1	46.6	49.1
LLaVA-1.5-7B	62.1	41.3	14.6	74.0	61.8	23.5	96.0	82.3	96.0	90.7
LLaMA-3-V-8B	60.0	36.4	17.5	70.1	43.6	20.6	97.8	83.2	99.0	90.7
Phi-3-V-3.8B	58.0	36.4	15.5	71.8	55.4	12.8	97.3	85.2	98.0	91.1
InstructBLIP-Vicuna-7B	37.6	25.7	0.0	29.9	15.2	0.0	55.0	51.4	47.8	50.2
Otter-Image-MPT7B	24.5	2.4	0.0	24.8	3.0	0.0	51.9	50.0	54.1	51.9
Random chance	25.0	6.3	0.4	25.0	6.3	0.4	50.0	50.0	50.0	50.0