

INTRIGUING PROPERTIES OF LARGE LANGUAGE AND VISION MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Recently, large language and vision models (LLVMs) have received significant attention and development efforts due to their remarkable generalization performance across a wide range of tasks requiring perception and cognitive abilities. A key factor behind their success is their simple architecture, which consists of a vision encoder, a projector, and a large language model (LLM). Despite their achievements in advanced reasoning tasks, their performance on fundamental perception-related tasks (e.g., MMVP) remains surprisingly low. This discrepancy raises the question of how LLVMs truly perceive images and exploit the advantages of the vision encoder. To address this, we systematically investigate this question regarding several aspects: *permutation invariance*, *robustness*, *math reasoning*, *alignment preserving* and *importance*, by evaluating the most common LLVM’s families (i.e., LLaVA) across 10 evaluation benchmarks. Our extensive experiments reveal several intriguing properties of current LLVMs: (1) they internally process the image in a global manner, even when the order of visual patch sequences is randomly permuted; (2) they are sometimes able to solve math problems without fully perceiving detailed numerical information; (3) the cross-modal alignment is overfitted to complex reasoning tasks, thereby, causing them to lose some of the original perceptual capabilities of their vision encoder; (4) the representation space in the lower layers ($< 25\%$) plays a crucial role in determining performance and enhancing visual understanding. Lastly, based on the above observations, we suggest potential future directions for building better LLVMs and constructing more challenging evaluation benchmarks.

1 INTRODUCTION

Large Language and Vision Models (LLVMs)(Liu et al., 2024c;b; Lee et al., 2024c;b) have demonstrated remarkable generalization capabilities across a wide variety of tasks, including coding and mathematics, showcasing their potential for practical applications. These impressive advancements have been achieved through a straightforward yet effective architecture based on the concept of *model-stitching*(Lenc & Vedaldi, 2015; Bansal et al., 2021). This approach integrates a pre-trained vision encoder (Radford et al., 2021) with a pre-trained large language model (LLM) (Touvron et al., 2023; Zheng et al., 2023b) via a simple cross-modal alignment module. This method significantly benefits from the power of well-established pre-trained representations. Consequently, this structure has become the *de facto* standard in the field, extending into other modality domains such as video, audio, and unified modalities (Xie et al., 2024; Erfei Cui, 2024).

Despite their significant generalization performance, recent studies have revealed several interesting phenomena regarding LLVMs. For instance, they struggle with tasks that are easy for humans to perceive (e.g., MMVP (Tong et al., 2024), BLINK (Fu et al., 2024)) and have limited understanding of domain-specific images (Zhai et al., 2024; Verma et al., 2024). In contrast to the computer vision domain, where demystifying the properties of vision encoders has been more thoroughly explored (Naseer et al., 2021; Kim et al., 2023; Vishniakov et al., 2023), the underlying properties of LLVMs are still largely under-explored. Therefore, in this work, we scrutinize the current *de facto* structure of LLVMs under various partial conditions, such as *permutation*, *occlusion*, and *synthetic images*.

In this paper, we systematically conduct a comprehensive evaluation of widely used LLVM families, specifically the LLaVA-series, across 10 diverse benchmarks, which encompass tasks such as math, chart, and basic perception. Our extensive experiments reveal several intriguing properties of current LLVMs, which we summarize as follows:

- In LLVMs, the visual patch tokens processed through the projector exhibit varying magnitudes of localized visual information. Remarkably, even when the order of these patch sequences is randomly shuffled before being fed into the LLM, the performance does not significantly degrade. For instance, in the case of LLaVA 1.5 (Li et al., 2024c), the average performance drop across 10 benchmarks is 0.19 ($< 1\%$), indicating that LLVMs exhibit permutation-invariant properties.
- LLVMs effectively handle tasks when given synthetic versions of the MathVista (Lu et al., 2023) dataset, with only a small performance decline (1.8% for LLaVA 1.5). Furthermore, we discovered that, in certain scenarios, LLVMs can solve problems even without access to the full image, including detailed numerical and chart elements.
- Following alignment and visual instruction tuning, LLVMs fail to preserve their initial perceptual capacities, with up to a 20% drop in image classification tasks (e.g., CIFAR-100 (Krizhevsky et al., 2009)), a phenomenon known as catastrophic forgetting (Zhai et al., 2024). Furthermore, they struggle to understand shared-world concepts within the representation space, according to the platonic representation hypothesis (Huh et al., 2024).
- Our analysis of model behavior reveals that LLVMs tend to concentrate more on the central region of the image. Furthermore, the lower layers in LLVM architectures are crucial for better generalization. In these layers (i.e., the bottom 20% of the LLM layers), the model primarily processes visual information, while the higher layers focus on interpreting the text.

In addition to our findings, we present and discuss several points regarding LLMs and evaluation benchmarks. Specifically, we highlight the need to develop more interactive and complex evaluation benchmarks to mitigate selection bias Zheng et al. (2023a) and improve applicability to real-world scenarios. Furthermore, when developing new LLMs, it is crucial to preserve cross-modal alignment. We hope that our findings will assist other ML researchers and engineers in building a new paradigm for LLMs.

2 RELATED WORKS

Large Language and Vision Models. Recent advancements in LLVMs have predominantly adopted simplistic yet highly effective architectures, notably through the model-stitching concept. Numerous prior studies have introduced various design modifications to bridge the performance gap with closed-source LLVMs (OpenAI, 2023; Anthropic, 2024). These efforts include focusing intently on high-resolution processing (Li et al., 2024e; Liu et al., 2024b; Shi et al., 2024), implementing locality-enhanced projectors (Cha et al., 2024), and incorporating knowledge embeddings (Lee et al., 2024c), layer traversal technique (Lee et al., 2024b) and leveraging a diverse array of vision encoders (Lu et al., 2024; Tong et al., 2024) have also been explored. Additionally, integrating external, task-specific computer vision modules (Lee et al., 2024d;e; Jiao et al., 2024; Lai et al., 2024) and incorporating different modalities — including video and audio (Wang et al., 2024; Li et al., 2024b; Erfei Cui, 2024; Xie et al., 2024) — have expanded the models’ capabilities. Moreover, enabling the handling of interleaved input formats (Li et al., 2024c; Xue et al., 2024) has further broadened the versatility of these models. While these models have been developed based on a simplistic structure, *model-stitching*, the effectiveness of this architecture remains under-explored.

Investigating Intriguing Properties of LLVMs. Alongside these advancements, recent studies have investigated and uncovered several crucial properties of current LLVMs. For instance, some studies have rigorously evaluated LLVMs on basic perception tasks that are trivially easy for humans by introducing “blind” pairs of image datasets (Tong et al., 2024; Fu et al., 2024; Rahmanzadehgervi et al., 2024). Other studies have explored cross-modal alignment by focusing on domain-specific visual capabilities (Verma et al., 2024) and examining the alignment of representation spaces across modalities between independently pre-trained LLMs and vision encoders (Li et al., 2024d; Huh

et al., 2024). Zhai et al. (2024) examine the phenomenon of catastrophic forgetting in LLMs within the context of image classification tasks. Additional studies (Zhou et al., 2023b; Chen et al., 2024b) analyze the persistent issue of object hallucination in LLMs. Moreover, research has explored spatial reasoning capabilities (Kamath et al., 2023). While vision encoders (e.g., ViT (Dosovitskiy, 2020), DeiT (Touvron et al., 2021)) in the computer vision field have been rigorously examined across a wide range of image settings, the study of these intriguing properties in LLMs remains relatively under-explored. In this paper, we aim to address this by conducting an in-depth investigation into LLMs, examining their permutation invariance, robustness, alignment preservation, and importance in scenarios involving occluded and synthesized images.

3 DEMYSTIFYING INTRIGUING PROPERTIES OF LLMs

In this section, we explore the intriguing properties of current LLMs that have *de facto* structure of *modal-stitching* in terms of various aspects: permutation invariance, robustness to occlusion, synthetic data, alignment preserving, and importance.

3.1 BACKGROUND

Overview of LLM. Current LLMs \mathcal{M} have widely adopted the *model-stitching* architecture, which consists of three main components: a pre-trained vision encoder f_v , a projector f_p , and a pre-trained LLM f_L . The overall model is represented as $\mathcal{M} = f_L \circ f_p \circ f_v$. The vision encoder f_v converts the input image $I \in \mathbb{R}^{3 \times H \times W}$ into visual features $\mathcal{F}_v \in \mathbb{R}^{N \times d_v}$, where $N = HW/P^2$ is the number of visual features, P is the patch size, and d_v is the dimension of the vision encoder’s output. The projector f_p transforms these visual features \mathcal{F}_v into visual patch tokens $\mathbf{X}_V \in \mathbb{R}^{N \times d_l}$ in the representation space of the LLM, where d_l is the embedding dimension of the LLM. This mapping allows the LLM to perceive and conceptually understand the given image. The LLM f_L produces an appropriate response $\mathbf{Y} = \{y_i\}_{i=1}^{L_Y}$ in an autoregressive manner, given both the visual patch tokens \mathbf{X}_V and the text tokens $\mathbf{X}_T \in \mathbb{R}^{L_T \times d_l}$, where L_T denotes the length of the input text sequence, and L_Y is the length of the output sequence. The probability of generating the response is given by:

$$p(\mathbf{Y} | \mathbf{X}_V, \mathbf{X}_T) = \prod_{i=1}^{L_Y} p(y_i | \mathbf{X}_V, \mathbf{X}_T, y_{<i}) \quad (1)$$

3.2 EVALUATION SETUP

Evaluation Benchmarks. To ensure a comprehensive and rigorous evaluation, we employ 10 standard and widely adopted benchmarks: MMVP (Tong et al., 2024), Q-Bench (Wu et al., 2023), MME (Fu et al., 2023), MMStar (Chen et al., 2024a), MM-Vet (Yu et al., 2023), LLaVA-W (Liu et al., 2024c), MathVista (Lu et al., 2023), SQA-IMG (Lu et al., 2022a), ChartQA (Masry et al., 2022), and AI2D (Kembhavi et al., 2016). Detailed descriptions of each dataset are provided in Appendix J.

Evaluation Models. Recently, a large number of LLM models have been actively introduced, owing to their remarkable flexibility and versatility across multiple domains. Consequently, it is challenging and inefficient to conduct holistic evaluations on all LLMs. Therefore, we select most standard LLMs: LLaVA-1.5-7B (Li et al., 2024c), LLaVA-NeXT-7B (Liu et al., 2024b), and LLaVA-OneVision-8B (Li et al., 2024b). For our customized experiments, before evaluating LLMs under diverse settings (e.g., occlusion), we first attempt to reproduce the baseline performance of LLMs on 10 evaluation benchmarks. To do this, we implement our customized evaluation toolkits by referring to the code of UniBench¹ (Al-Tahan et al., 2024). Detailed descriptions of each model are provided in Appendix K.

¹<https://github.com/facebookresearch/unibench>

LLVMs	MMVP	Q-Bench	MME	MMStar	MM-Vet	LLaVA ^W	MathVista	SQA ^f	ChartQA	A12D	Avg. Δ
LLaVA-1.5	34.67	59.73	1850.07	34.20	31.50	67.50	24.70	65.59	16.92	53.34	
+ Perm.	36.00	59.60	1874.60	33.33	30.40	66.20	21.20	65.44	14.08	52.69	▼ 0.59
	(▲ 1.33)	(▼ 0.13)	(▲ 24.53)	(▼ 0.87)	(▼ 1.10)	(▼ 1.30)	(▼ 3.50)	(▼ 0.15)	(▼ 2.84)	(▼ 0.65)	
LLaVA-NeXT	36.67	63.55	1874.42	37.80	43.50	75.50	32.00	62.12	66.06	64.02	
+ Perm.	37.33	62.54	1890.19	36.87	43.40	75.80	21.70	62.12	34.55	64.02	▼ 2.71
	(▲ 0.67)	(▼ 1.00)	(▲ 15.78)	(▼ 0.93)	(▼ 0.10)	(▲ 0.30)	(▼ 10.30)	(▼ 0.00)	(▼ 31.51)	(▼ 0.00)	
LLaVA-OneVision	60.67	77.26	1982.5	59.87	57.80	87.40	61.80	94.00	93.52	81.25	
+ Perm.	59.33	76.99	1964.3	54.93	47.60	82.30	53.50	89.24	58.26	75.58	▼ 9.40
	(▼ 1.33)	(▼ 0.27)	(▼ 18.2)	(▼ 4.93)	(▼ 10.20)	(▼ 5.10)	(▼ 8.30)	(▼ 4.76)	(▼ 35.26)	(▼ 5.67)	

Table 1: Results of drop ratio (Δ) when random permutation is applied. We run five experiments.

3.3 DO LLVMS PERCEIVE IMAGES GLOBALLY?

Current LLMs commonly adopt ViT (Dosovitskiy, 2020)-based vision encoders, such as CLIP ViT (Radford et al., 2021) and SigLIP (Zhai et al., 2023), making their image perception dependent on these encoders. Specifically, ViT is designed to learn interactions across all image patches, providing properties (Naseer et al., 2021; Vishniakov et al., 2023) such as *permutation invariance* and *robustness to occlusion*. This raises the question of whether these ViT properties might transfer to current LLM models.

Each visual patch token encapsulates localized visual information.

We first investigate whether each visual patch token X_V from the projector f_P captures a localized understanding of the patch area corresponding to its position in the image. Specifically, given an image I , the projector outputs N visual patch tokens (e.g., $N = 576$ for LLaVA-1.5-7B). We then select a single token (removing all others) and feed it into the LLM f_L . To quantify this, we define the patch information loss (PIL) as the ratio of the performance drop to the original performance. However, performing computations on each individual visual token is computationally intensive, especially for models such as LLaVA-1.5-7B that process 576 visual tokens arranged in a 24×24 grid of patches. To accelerate computation and reduce complexity, we aggregate the original N visual tokens into M tokens, where $M < N$, by grouping neighboring tokens. As shown in Figure 1, the group-wise visual tokens in the LLaVA-1.5-7B model demonstrate varying levels of performance on the MMStar (Chen et al., 2024a) and MME (Fu et al., 2023). Darker regions indicate areas where the model retains more localized information for those specific groups.

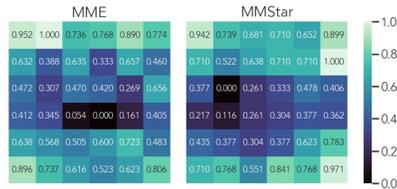


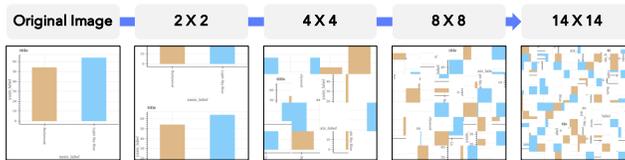
Figure 1: We demonstrate the extent to which group-wise visual tokens capture region-specific information (PIL) for LLaVA-1.5-7B on the MMStar (Chen et al., 2024a) and MME (Fu et al., 2023). Darker regions indicate areas where the model retains more localized information for those specific groups.

LLVMs are permutation invariant in terms of visual patch tokens, depending on the benchmark.

From our above results, we empirically verify that each visual patch token from the projector contains localized visual information. Here, we aim to understand how LLVMs systematically process and perceive images based on these visual patch tokens. Given that LLVMs generate answers in an autoregressive manner, we investigate whether they exhibit *order bias* regarding visual patch tokens. To study this, we strongly hypothesize that if LLVMs have *permutation variance*, the performance drop (Δ) will be significant when a random permutation is applied to the visual patch tokens X_V .

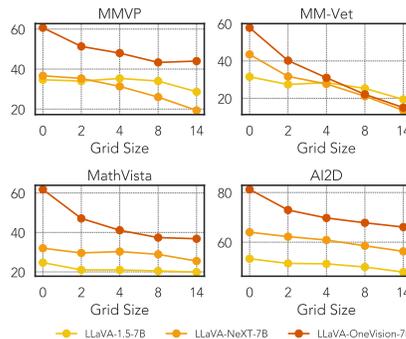
As shown in Table 1, the overall performance across most benchmarks declines when the visual patch tokens are randomly shuffled. However, the performance gap between the original and the shuffled (Perm.) versions is not substantial, remaining within a 0–2% range, for LLaVA-1.5 and LLaVA-NeXT. Considering that LLaVA-1.5 uses 576 visual tokens, this is an intriguing observation. It suggests that current LLVMs interpret images in a *global* manner, despite each visual patch token containing localized information (see Figure 1), and even though they process both images and text autoregressively. In the case of LLaVA-OneVision which has many visual tokens (729), the avg. performance drop (Δ) is non-trivial. Upon closer analysis, we find that the “permutation invariance”

216 Figure 3: We present examples of
 217 shuffled images with different grid
 218 sizes (2, 4, 8, 14) derived from a
 219 MathVista dataset image. As the grid
 220 size increases, the chart image be-
 221 comes more artistically styled.



222
 223
 224 property is both benchmark-dependent and capability-dependent. Specifically, in perception-related
 225 benchmarks, such as MMVP and Q-Bench, the performance drop is minimal. In fact, for LLaVA-1.5
 226 and LLaVA-NeXT, performance even slightly improves in some cases. On the other hand, in text-rich
 227 benchmarks requiring reasoning capabilities (e.g., MathVista and ChartQA), the performance drops
 228 significantly. These benchmarks demand an understanding of detailed numerical information and
 229 highly structured geometric graphs, where preserving the spatial structure of visual patch tokens is
 230 crucial. We hypothesize that this global interpretation may result from recent LLMs being trained
 231 via backpropagation, with the loss signal primarily derived from the text output of the Assistant:
 232 side. Based on these experiments, we argue that while LLMs are trained with an autoregressive
 233 objective, they internally handle images globally. This observation offers a possible explanation
 234 for the success of pixel shuffling (Chen et al., 2024c) in achieving both strong performance and
 efficiency.

235 **LLMs are sensitive to spatial structures.** In-
 236 stead of treating visual patch tokens as permutation
 237 invariants, we explore how LLMs behave when the
 238 sequence of image patches is permuted. To exam-
 239 ine the sensitivity to spatial structure, we randomly
 240 shuffle image patches at varying grid sizes (2, 4,
 241 8, 14), as shown in Figure 2. In our experiments,
 242 we observe that LLaVA-OneVision is sensitive to
 243 spatial structures on the MathVista (Lu et al., 2023)
 244 and AI2D (Kembhavi et al., 2016) tasks, despite the
 245 ViT learning all interactions between image patches.
 246 This result contrasts with previous study (Naseer
 247 et al., 2021) suggesting that ViT-based vision en-
 248 coders exhibit high permutation invariance to patch
 249 positions than CNN counterparts. We posit that on
 250 the MMVP Tong et al. (2024) dataset, which in-
 251 volves perception task, LLaVA-OneVision would
 252 also show permutation invariance to randomly shuffled
 253 patches, similar to existing work (Naseer
 254 et al., 2021) analyzing the ImageNet (Deng et al.,
 255 2009) val. dataset. However, unlike ImageNet, the
 256 MathVista and AI2D datasets contain more struc-
 257 turally complex images (e.g., charts, code screen-
 258 shots) that are highly sensitive to spatial structure,
 259 as the original numerical understanding is signifi-
 260 cantly disrupted. **Shuffling image patches in such cases disrupts geometric relationships or relative
 magnitudes in plots or charts, making accurate interpretation of these images significantly more
 challenging.** Interestingly, both LLaVA-1.5 and LLaVA-NeXT exhibit insensitivity to spatial structure,
 particularly in the MathVista dataset, where performance drops were minimal. These results
 suggest the need for further investigation, which we address in the following sections.



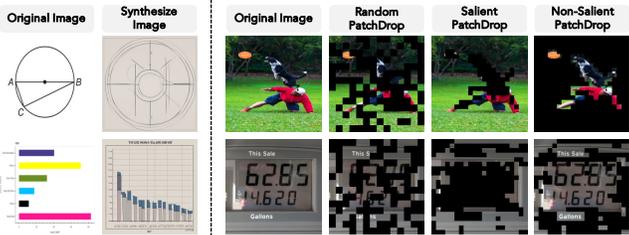
261 Figure 2: We present the performance across
 262 different grid sizes (2, 4, 8, 14) on the
 263 MMVP, MM-Vet, MathVista, and AI2D
 264 datasets, using three models: LLaVA-1.5,
 265 LLaVA-NeXT, and LLaVA-OneVision.

261 3.4 DO LLMs PERCEIVE NUMERICAL INFORMATION WELL?

262
 263 Here, we study whether LLMs truly perceive text-rich images (e.g., charts, geometric shapes) that
 264 contain highly detailed numerical and shape information. To do this, we construct synthetic datasets
 265 for MMVP (Tong et al., 2024) and MathVista (Lu et al., 2023). Specifically, we first generate an
 266 image description of a given image using LLaVA-OneVision (Li et al., 2024b) with the prompt:
 267 “Please generate a caption of this image.”. Next, we generate an image corresponding to the image
 268 description leveraging the SDXL-Lightning (Lin et al., 2024) model, ensuring both quality and
 269 efficiency. As a result, we get synthesized version (Syn.) to the original version (Org.), illustrated
 in Figure 5.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

Figure 5: We present examples of images (left) synthesized by SDXL-Lightning and (right) occluded using three methods: Random, Salient, and Non-Salient. The original images are from the MathVista and MME datasets. Occluded areas are marked in black to indicate zero pixel values.



In some cases, LLMs can solve problems without seeing the image. Table 2 presents the performance comparison on both original and synthesized datasets. For comparison, we evaluate the knowledge-embedded-specific LLM, Meteor 7B (Lee et al., 2024c). Overall, compared to the basic perception task (i.e., MMVP), the performance drop in MathVista is not significantly larger across four LLMs. Given that the generated images show distorted chart and function shapes, with detailed numerical and formula information missing, as shown in Figure 5 (CLIP-I scores lower than in MMVP), it is surprising that LLMs are still able to solve math problems requiring advanced cognitive reasoning, without key information. This observation leads us to more in-depth analysis on MathVista dataset. We analyze how LLMs solve math problems using synthesized images. In instances where they answer correctly, LLMs frequently choose “No” for MCQs and tend to generate “1” for free-form responses. A deeper analysis reveals that many of these questions ask “What is the smallest value?”, causing the models to select “1” using commonsense reasoning, without needing to interpret the image itself. Table 2 shows how often the models produce “1,” with a noticeable drop in frequency for LLaVA-OneVision and Meteor models. This suggests that these models, likely due to extensive training with million-scale datasets, struggle with “smallest value” questions when images are unclear, demonstrating their ability to interpret images effectively.

Scaling up visual instruction tuning datasets improves text-only mathematical reasoning. Here, we explore whether enhancing math reasoning in a visual context can improve standard text-only math reasoning. We evaluate four LLMs on the GSM8K (Cobbe et al., 2021) dataset in an 8-shot setting using Chain-of-Thought (CoT) prompting (Wei et al., 2022). As shown in Figure 4, we observe that models performing well in visual math contexts also achieve strong performance on GSM8K. Moreover, as the scale of the dataset used for training increases, so does model performance. This suggests that using high-quality, large-scale datasets (e.g., rationale-style datasets, as used in Meteor) is beneficial, and that there is compatibility between visual math and text-only math reasoning, aligning with the data-centric AI perspective (Xu et al., 2023; Zhou et al., 2024).

3.5 ARE LLMs ROBUST TO OCCLUSIONS?

Existing studies (Naseer et al., 2021; Vishniakov et al., 2023) have demonstrated that ViT models exhibit a remarkable degree of robustness to occlusions, such as patch dropping, than CNN counterparts. Since most LLMs utilize CLIP ViT-L as their vision encoder, we aim to explore

LLMs	MMVP				MathVista		
	Text-rich?	Scale	Orig.	Syn.	Orig.	Syn.	Freq.
CLIP-I	-	-	-	0.84	-	0.61	-
LLaVA-1.5	✗	158K	34.7	20.7 (▼14.0)	24.7	22.9 (▼1.8)	81.0
LLaVA-NeXT	✓	760K	36.7	16.7 (▼20.0)	32.0	27.7 (▼4.3)	50.0
Meteor	✓	1.1M	51.3	35.3 (▼16.0)	52.1	31.4 (▼20.7)	9.5
LLaVA-OneVision	✓	3.1M	60.7	37.3 (▼23.3)	61.8	37.0 (▼24.8)	12.0

Table 2: We present the performance on the synthesized versions of the MMVP (Tong et al., 2024) and MathVista (Lu et al., 2023) datasets across the models LLaVA-1.5, LLaVA-NeXT, Meteor, and LLaVA-OneVision. Additionally, we provide the scale of the visual instruction training datasets used by each model and specify whether chart, math, and diagram datasets were included. CLIP-I indicates the image similarity using CLIP-ViT-L/14. Freq. denotes the frequency with which the model generates the answer “1” in free-form question types in Syn. cases.

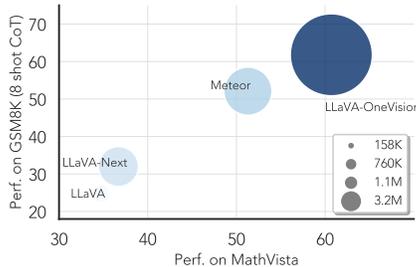


Figure 4: We present performance on the GSM8K dataset using 8-shot Chain-of-Thought prompting. Additionally, we demonstrate that scaling up the instruction-tuning dataset enables LLMs to solve text-only math reasoning problems more effectively.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

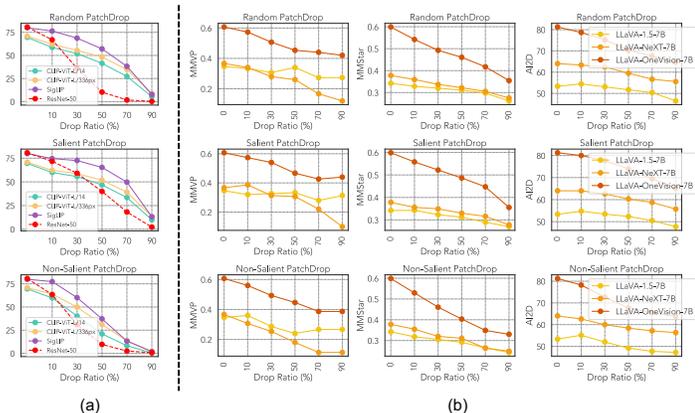


Figure 6: We present robustness performance under occlusion conditions. (a) ViT variant vision encoders demonstrate greater robustness to occlusion compared to ResNet-50. (b) LLMs also show robustness to occlusion, benefiting from the use of ViT encoders.

whether this robustness transfers to LLMs in scenarios involving occluded images. Following the simple masking method presented in prior work (Naseer et al., 2021), we manipulate the input image $I = \{x_i\}_{i=1}^N$, where N represents the number of patches. Specifically, we mask out N' patches (where $N' < N$) by setting the pixel values of these patches to zero, creating an occluded image, denoted as \tilde{I} . We then apply three distinct occlusion methods to the image I : (1) **Random PatchDrop**: A subset of N' patches is randomly selected and dropped from the image, effectively simulating random occlusion; (2) **Salient PatchDrop**: We strategically select and drop salient patches by leveraging the self-supervised ViT model `dino-small` (Caron et al., 2021); (3) **Non-Salient PatchDrop**: In this case, we drop non-salient, background patches, retaining only the salient information. This method also utilizes `dino-small`, following a similar approach to the **Salient PatchDrop** but focusing on removing the background regions. Figure 5 presents example images with different occlusion methods applied.

LLMs are robust against occlusion. Before evaluating LLMs on occluded images, we first verify whether ViT-based encoders are more robust than their CNN counterparts in this scenario. To do this, we assess several ViT variants and ResNet-50 (He et al., 2016) on occluded ImageNet (Krizhevsky et al., 2009) images, applying the same masking process as mentioned above. As shown in Figure 6 (left), compared to ResNet-50, ViT variants demonstrate greater robustness in occlusion scenarios, consistent with findings from the prior study (Naseer et al., 2021). Due to this robustness, LLMs also exhibit relatively strong performance under occlusion. This result is surprising given that in the AI2D dataset, which contains text-rich diagram images, 50%–70% of the patches are missing, yet LLMs can still provide correct answers to some extent. This may be because AI2D involves selecting one answer from multiple options, suggesting the possibility of selection bias (Zheng et al., 2023a), a significant issue that we leave for future work.

3.6 DO LLMs PRESERVE CROSS-MODAL ALIGNMENT?

In the *de facto* structure of LLMs, a projector f_P enables LLMs to perceive and understand images by transforming visual representations into the LLM’s representation space. While a recent work (McKinzie et al., 2024) suggests that the type of projector has minimal impact on performance, other studies (Zhai et al., 2024; Verma et al., 2024) have argued that the projector have limitations in preserving cross-modal understanding and issues such as catastrophic forgetting. In this work, we investigate (1) how effectively a trained projector maintains its *visual recognition* capability relative to the LLM’s original vision encoders (e.g., CLIP-ViT-14 for LLaVA-NeXT), and (2) how well a trained projector preserves cross-modal alignment, based on the *platonic representation hypothesis* (Huh et al., 2024), compared to representation expressivity without alignment learning.

LLMs struggle to preserve the original visual understanding capability. Ideally, after alignment and visual instruction tuning, LLMs should retain the visual perception abilities of their original vision encoders, allowing them to understand and classify images effectively. To assess this, we evaluate LLMs on zero-shot image classification tasks using widely adopted datasets such as Caltech100 (Higgins et al., 2017), Food101 (Bossard et al., 2014), CIFAR-100 (Krizhevsky

LLVMs	Caltech101	CIFAR-100	Food101	Pets	Country211	EuroSAT	AirCRAFT	Avg.
CLIP ViT 336	84.50 (0.52)	75.10 (1.74)	93.72 (0.61)	93.48 (0.63)	31.14 (0.84)	58.90 (0.84)	34.08 (0.77)	67.27 (0.85)
LLaVA-1.5	43.76 (2.69) (▼ 40.74)	48.36 (2.47) (▼ 26.74)	57.22 (0.61) (▼ 36.5)	73.22 (0.37) (▼ 20.26)	12.20 (0.47) (▼ 18.94)	11.72 (0.34) (▼ 47.18)	17.00 (0.72) (▼ 17.08)	37.64 (1.10) (▼ 29.63)
CLIP ViT 14	84.52 (0.56)	75.86 (1.06)	92.78 (0.41)	93.08 (0.30)	28.68 (1.44)	58.46 (0.80)	32.98 (1.02)	66.62 (0.80)
LLaVA-NeXT	56.68 (2.29) (▼ 27.84)	45.36 (1.02) (▼ 30.5)	53.14 (1.00) (▼ 39.64)	75.06 (0.93) (▼ 18.02)	12.94 (0.35) (▼ 15.74)	8.34 (0.59) (▼ 50.12)	12.66 (0.46) (▼ 20.32)	37.74 (▼ 28.88)

Table 3: We report the Top-1 accuracy (%) with standard deviation (in parentheses) of LLVMs and their corresponding vision encoder models on 1K subsampled datasets from Caltech100, CIFAR-100, Food101, Pets, Country211, EuroSAT, and AirCRAFT. We run five experiments.

et al., 2009), Pets (Parkhi et al., 2012), Country211², EuroSAT (Helber et al., 2019), and Air-Craft (Maji et al., 2013). Following the method in previous work (Zhai et al., 2024), we use ChatGPT (gpt-3.5-turbo) (OpenAI, 2023) to extract a single label with the use of prompt: *Is this prediction correct?*. As shown in Table 3, the performance of LLVMs significantly degrades across all datasets compared to their vision encoders, suggesting that LLVMs do not fully retain the perception capabilities of their original vision encoders. This may be due to: (1) LLVMs being trained to solve complex tasks (e.g., chart or math reasoning) with the use of instruction, which may cause them to lose basic perception abilities (e.g., recognizing simple objects), a phenomenon known as *catastrophic forgetting* (Zhai et al., 2024)³, and (2) the vision encoder’s relatively small parameter size (307M for CLIP ViT-L/336px) compared to the LLM (7B for Vicuna-1.5), which could result in a loss of visual perception capability during projection, as the more powerful LLM dominates.

LLVMs lose the ability to understand and interpret shared world concepts.

Beyond visual recognition capabilities, we analyze cross-modal alignment based on the *platonic representation hypothesis* (Huh et al., 2024), which argues that neural networks, despite being trained on different objectives, data, and modalities, should converge to a shared statistical model of reality in their representation spaces. To measure representation similarity between two modalities, the original authors of this hypothesis use mutual nearest-neighbor alignment metrics, a type of kernel-alignment metric. In our work, we assess how much alignment is lost after visual instruction tuning by applying this metric within the context of the platonic representation hypothesis. We evaluate 10 LLMs and measure alignment between these LLMs and vision encoders (LLVMs) using the DOCCI (Onoe et al., 2024) dataset which contains long image descriptions requiring localized scene understanding. As shown in Figure 7, after visual instruction tuning, both LLaVA-1.5 and LLaVA-NeXT show degraded alignment performance with respect to representations compared to their original vision encoder. This suggests doubts about the actual role of the projector in causing the degradation in alignment preservation. From this observation, we speculate that current LLVMs are trained on a variety of datasets to achieve generalization (i.e., multi-task learning). However, during visual instruction tuning, the models might overemphasize capabilities requiring complex cognition while potentially reducing representations related to other tasks, such as localized scene understanding (i.e., DOCCI). This results in a lower alignment score and catastrophic forgetting, as shown in Table 3. For future work, one potential direction is to develop a localized enhanced alignment module similar to HoneyBee (Cha et al., 2024).

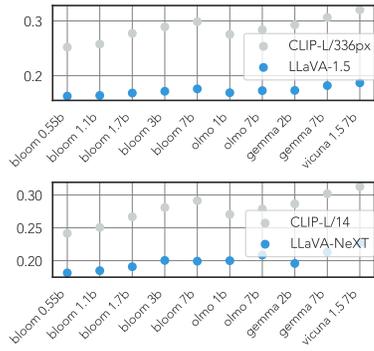


Figure 7: We present how alignment preservation changes (CLIP → LLaVA) in the representation space across various LLM families, BLOOM (Le Scao et al., 2023), OLMo (Groeneveld et al., 2024), Gemma (Team et al., 2024), Vicuna (Chiang et al., 2023), with different parameter sizes on the DOCCI dataset.

²<https://github.com/openai/CLIP/blob/main/data/country211.md>
³On the CLEVR/Count (Johnson et al., 2017) dataset, we observed a 16.6% performance improvement in the LLaVA-NeXT model compared to the previous vision encoder (i.e., CLIP-ViT-L/14.)

3.7 MODEL BEHAVIOR: WHICH MODALITY AND LAYERS ARE MORE IMPORTANT?

Here, we conduct an in-depth analysis of model behavior to assess the importance of either a layer or a visual token when performing downstream tasks. We hypothesize that if adding arbitrary noise to a specific component—either a layer block or a visual token—results in a significant drop in model performance, then that component is crucial and actively involved in the model’s reasoning process. To quantify this, we define an *importance* score (\mathcal{I}) inspired by the concept of “sharpness of minima” (Keskar et al., 2016; Lee et al., 2024f). This concept aims to identify flat minima, which promote stable training and better generalization, by measuring the sensitivity of the training function around a local minimum. In our work, we adapt this concept for the inference stage.

Definition 3.1 (Importance Score). Let $x_t \in \mathbb{R}^d$ be the d -dimensional input embedding for a target subject t . For x_t , we define the constraint candidate set \mathcal{C}_t for t as:

$$\mathcal{C}_t = \{z_t \in \mathbb{R}^d : -\epsilon + |x_t| \leq z_t \leq \epsilon + |x_t|\}, \quad \epsilon \sim \text{Uniform}(-1, 1), \tag{2}$$

where z_t is a noise vector. The importance score \mathcal{I} for target t is then defined as:

$$\mathcal{I}_t := \frac{f(x_t) - \max_{z_t \in \mathcal{C}_t} f(x_t + z_t)}{f(x_t)} \times 100. \tag{3}$$

Note that while the concept of “sharpness of minima” was originally used to find flat minima during training by defining a square-bound constraint set, this is feasible because the model is trained via stochastic gradient descent (SGD), which indirectly allows the evaluation of all noise values z in the constraint set C . However, since our experiment focuses on downstream task performance during inference, we adapt this concept by sampling K noise vectors, $\{z_t^1, z_t^2, \dots, z_t^K\} \sim \mathcal{C}_t$, with different random seeds. For computational efficiency, we set $K = 10$.

LLVMs strongly focus on the center of the image. To assess the extent to which LLVM utilizes visual token information, we add a noise vector z_t to each visual token information based on the importance score (\mathcal{I}). However, performing computations on each individual visual token is computationally intensive, especially for models such as LLaVA-1.5-7B that process 576 visual tokens arranged in a 24×24 grid of patches. To accelerate computation and reduce complexity, we adopt the same process as 3.3. As shown in Figure 8, LLaVA-1.5-7B places strong emphasis on the central part of the images in the MM-Vet dataset, while the edge regions have minimal influence on the final performance compared to the central visual tokens. This result suggests that not all visual tokens are necessary, aligning with recent works (Cha et al., 2024; Xue et al., 2024) that reduce redundant visual tokens in the projector to enhance efficiency.

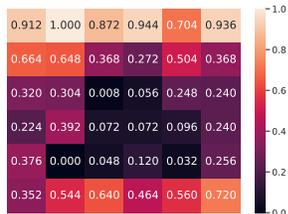


Figure 8: We report the degree of utility of group-wise visual tokens for LLaVA 1.5 7B on the MM-Vet dataset. Darker regions indicate that the LLVM relies heavily on information from those specific group parts.

Lower block layers in LLVMs are more important.

Figure 9 (left) shows that the lower layers (< 6) play a crucial role in handling the integrated capabilities required for tasks in the MMStar dataset. This suggests that these layers contain more beneficial representations for perception and cognition. This finding aligns with recent work on LLVMs, specifically TroL (Lee et al., 2024b), which introduces the concept of “layer traversal.” This technique revisits layer representations, resulting in a highly generalizable model despite its small size (1.8B parameters). In their paper (Figure 6), the traversal pattern is more pronounced in the lower layers, which is consistent with our findings. Therefore, we believe our results may provide insights into why traversing layers leads to improved generalization.

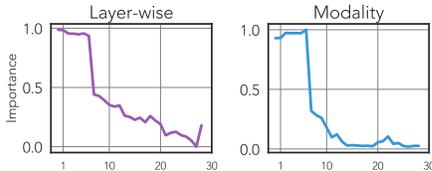


Figure 9: We present the results of (left) layer-wise importance and (right) modality importance within the layers.

Textual modality is more important than visual modality. In addition to layer-wise importance, we measure *modality importance* using the score $\frac{\mathcal{I}_T}{\mathcal{I}_I}$, which calculates the relative importance of textual and image modalities. Specifically, to obtain the image modality importance score \mathcal{I}_I , we feed

noise vectors to the positions corresponding to image tokens (e.g., 576 tokens in LLaVA-1.5), and vice versa for text modality \mathcal{I}_T . As shown in Figure 9 (right), until the lower layers (< 8), the image modality is more important than the textual modality. However, as layers progress, we observe that the textual modality becomes increasingly important, likely because generating responses requires a stronger focus on text with the perspective of autoregressive modeling. This suggests that LLVMs initially focus on global image perception (section 3.3), and by the middle layers, they have processed the image and shifted toward solving complex user queries to generate coherent answers. Similarly, in TroL, layer traversal occurs more actively in the lower layers, which we interpret as the model attempting to better comprehend the image when it fails to do so in a single pass, enabling it to solve complex reasoning tasks more effectively. These findings highlight the value of strong visual perception, which may explain the success of models utilizing large visual tokens (Wang et al., 2024; Li et al., 2024b) or high-resolution image processing (Li et al., 2024b).

4 DISCUSSIONS

Building more interactive evaluation benchmarks. As mentioned in section 3.4, LLVM can effectively solve problems even without seeing the input image. However, current evaluation benchmarks are designed for single-turn interactions and lack applicability to real-world, interactive scenarios. For example, in standard OCR tasks, we typically assess whether the LLVM correctly transcribes text from an image. But consider a practical situation: you’re traveling in a foreign country and visiting a local restaurant. Translating the menu is challenging, and while an application with strong OCR capabilities would be helpful, this is only the first step. When ordering, the LLVM should not only recognize the menu items but also understand the user’s preferences — what they like or dislike — by incorporating knowledge of their persona (Lee et al., 2022). Therefore, future benchmarks should be more interactive and socially grounded (Zhou et al., 2023a), extending beyond multiple-choice, binary, or non-interactive free-form tasks. These benchmarks should involve multi-turn interactions and be based on the user’s preferences (Lee et al., 2024g) or persona in long-term social interactions (Jang et al., 2023; Lee et al., 2024h).

A new paradigm enhancing cross-modal alignment. Current LLVMs have widely adopted the *model-stitching* structure, which demonstrates impressive capabilities on tasks requiring higher-level cognitive reasoning. However, they exhibit significantly degraded performance in zero-shot image classification tasks (Table 3). Additionally, they cannot effectively preserve alignment (Figure 7) in terms of relatively simple perception when compared to text-rich images (e.g., charts, mathematical expressions). Recently, although recent studies (Li et al., 2024b; Lee et al., 2024c) has been extensively scaling up model sizes with larger datasets to achieve higher performance on increasingly difficult tasks — which we believe is the correct direction - we think it is necessary to deeply consider innovative model architectures (e.g., layer of traversal (Lee et al., 2024b), hidden dimension expansion (Lee et al., 2024a)) to enhance cross-modal alignment at least once. For example, in recent unified architectures (Xie et al., 2024; Erfei Cui, 2024), enabling LLMs to generate images is akin to how drawing can be substantially more challenging for humans than simply viewing a picture. This is because drawing requires a comprehensive and simultaneous understanding of complex world concepts such as relationships between objects, lighting, perspective, and more. Therefore, by projecting visual imagination abilities (Lu et al., 2022b; Lee et al., 2023) onto LLVMs to enable them to generate images, it might significantly help in better preserving cross-modal alignment.

5 CONCLUSION

In this paper, we systemically reveals intriguing properties of LLVMs with respect to *permutation invariance*, *robustness*, *alignment preserving*, and *importance* under various image settings such as occlusion and synthesized images. We hope these findings will assist academic researchers and ML developers in advancing the next frontier of LLMs by providing a foundational basis for future model design choices.

REFERENCES

Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, and Peter Anderson. Nocaps: Novel object captioning at scale. In

- 540 *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 8948–8957, 2019.
- 541
- 542 Haider Al-Tahan, Quentin Garrido, Randall Balestriero, Diane Bouchacourt, Caner Hazirbas, and
- 543 Mark Ibrahim. Unibench: Visual reasoning requires rethinking vision-language beyond scaling.
- 544 *arXiv preprint arXiv:2408.04810*, 2024.
- 545 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel
- 546 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language
- 547 model for few-shot learning. *Advances in neural information processing systems*, 35:23716–
- 548 23736, 2022.
- 549 Anthropic. The claude 3 model family: Opus, sonnet, haiku. <https://www.anthropic.com>,
- 550 2024. URL https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf.
- 551
- 552
- 553 Yamini Bansal, Preetum Nakkiran, and Boaz Barak. Revisiting model stitching to compare neural
- 554 representations. *Advances in neural information processing systems*, 34:225–236, 2021.
- 555 Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani,
- 556 and Sağnak Taşirlar. Introducing our multimodal models, 2023. URL <https://www.adept.ai/blog/fuyu-8b>.
- 557
- 558 Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative compo-
- 559 nents with random forests. In *Computer vision—ECCV 2014: 13th European conference, zurich,*
- 560 *Switzerland, September 6-12, 2014, proceedings, part VI 13*, pp. 446–461. Springer, 2014.
- 561
- 562 Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui
- 563 Chen, Zhi Chen, Pei Chu, et al. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*,
- 564 2024.
- 565 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and
- 566 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of*
- 567 *the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- 568
- 569 Junbum Cha, Wooyoung Kang, Jonghwan Mun, and Byungseok Roh. Honeybee: Locality-enhanced
- 570 projector for multimodal llm. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
- 571 *and Pattern Recognition*, pp. 13817–13827, 2024.
- 572 Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing
- 573 web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the*
- 574 *IEEE/CVF conference on computer vision and pattern recognition*, pp. 3558–3568, 2021.
- 575 Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua
- 576 Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint*
- 577 *arXiv:2311.12793*, 2023a.
- 578
- 579 Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi
- 580 Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language
- 581 models? *arXiv preprint arXiv:2403.20330*, 2024a.
- 582 Xuweiyi Chen, Ziqiao Ma, Xuejun Zhang, Sihan Xu, Shengyi Qian, Jianing Yang, David F
- 583 Fouhey, and Joyce Chai. Multi-object hallucination in vision-language models. *arXiv preprint*
- 584 *arXiv:2407.06192*, 2024b.
- 585 Yingyi Chen, Xi Shen, Yahui Liu, Qinghua Tao, and Johan AK Suykens. Jigsaw-vit: Learning
- 586 jigsaw puzzles in vision transformer. *Pattern Recognition Letters*, 166:53–60, 2023b.
- 587
- 588 Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong,
- 589 Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to com-
- 590 mercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024c.
- 591 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,
- 592 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna:
- 593 An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.

- 594 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
595 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reichihiro Nakano, et al. Training verifiers to
596 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
597
- 598 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-
599 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
600 pp. 248–255. Ieee, 2009.
- 601 Karan Desai, Gaurav Kaul, Zubin Aysola, and Justin Johnson. Redcaps: Web-curated image-text
602 data created by the people, for the people. *arXiv preprint arXiv:2111.11431*, 2021.
603
- 604 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning
605 of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- 606 Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale.
607 *arXiv preprint arXiv:2010.11929*, 2020.
- 608 et al. Erfei Cui. Sharegpt-4o: Comprehensive multimodal annotations with gpt-4o, 2024. <https://sharegpt4o.github.io/>.
609
610
- 611 Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu
612 Zheng, Ke Li, Xing Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal
613 large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- 614 Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A
615 Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but
616 not perceive. *arXiv preprint arXiv:2404.12390*, 2024.
- 617 Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord,
618 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating the
619 science of language models. *arXiv preprint arXiv:2402.00838*, 2024.
620
- 621 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
622 *preprint arXiv:2312.00752*, 2023.
- 623 Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey,
624 and Noah A Smith. Don’t stop pretraining: Adapt language models to domains and tasks. *arXiv*
625 *preprint arXiv:2004.10964*, 2020.
626
- 627 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
628 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
629 770–778, 2016.
- 630 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked au-
631 toencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer*
632 *vision and pattern recognition*, pp. 16000–16009, 2022.
- 633 Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset
634 and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected*
635 *Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.
636
- 637 Irina Higgins, Loic Matthey, Arka Pal, Christopher P Burgess, Xavier Glorot, Matthew M Botvinick,
638 Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a
639 constrained variational framework. *ICLR (Poster)*, 3, 2017.
- 640 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
641 and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint*
642 *arXiv:2106.09685*, 2021.
- 643 Minyoung Huh, Brian Cheung, Tongzhou Wang, and Phillip Isola. The platonic representation
644 hypothesis. *arXiv preprint arXiv:2405.07987*, 2024.
645
- 646 Neel Jain, Ping-yeh Chiang, Yuxin Wen, John Kirchenbauer, Hong-Min Chu, Gowthami Somepalli,
647 Brian R Bartoldson, Bhavya Kailkhura, Avi Schwarzschild, Aniruddha Saha, et al. Neftune:
Noisy embeddings improve instruction finetuning. *arXiv preprint arXiv:2310.05914*, 2023.

- 648 Jihyoung Jang, Minseong Boo, and Hyounghun Kim. Conversation chronicles: Towards
649 diverse temporal and relational dynamics in multi-session conversations. *arXiv preprint*
650 *arXiv:2310.13420*, 2023.
- 651 Qirui Jiao, Daoyuan Chen, Yilun Huang, Yaliang Li, and Ying Shen. Enhancing multimodal
652 large language models with vision detection models: An empirical study. *arXiv preprint*
653 *arXiv:2401.17981*, 2024.
- 654 Xisen Jin, Dejiao Zhang, Henghui Zhu, Wei Xiao, Shang-Wen Li, Xiaokai Wei, Andrew Arnold, and
655 Xiang Ren. Lifelong pretraining: Continually adapting language models to emerging corpora.
656 *arXiv preprint arXiv:2110.08534*, 2021.
- 657 Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and
658 Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual
659 reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
660 pp. 2901–2910, 2017.
- 661 Amita Kamath, Jack Hessel, and Kai-Wei Chang. What’s” up” with vision-language models? in-
662 vestigating their struggle with spatial reasoning. *arXiv preprint arXiv:2310.19785*, 2023.
- 663 Zixuan Ke, Bing Liu, and Xingchang Huang. Continual learning of a mixed sequence of similar and
664 dissimilar tasks. *Advances in neural information processing systems*, 33:18493–18504, 2020.
- 665 Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali
666 Farhadi. A diagram is worth a dozen images. In *Computer Vision—ECCV 2016: 14th Euro-
667 pean Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*,
668 pp. 235–251. Springer, 2016.
- 669 Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Pe-
670 ter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv*
671 *preprint arXiv:1609.04836*, 2016.
- 672 Junho Kim, Byung-Kwan Lee, and Yong Man Ro. Demystifying causal features on adversarial ex-
673 amples and causal inoculation for robust network by adversarial instrumental variable regression.
674 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
675 12302–12312, 2023.
- 676 James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A
677 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcom-
678 ing catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*,
679 114(13):3521–3526, 2017.
- 680 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
681 2009.
- 682 Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. Lisa: Rea-
683 soning segmentation via large language model. In *Proceedings of the IEEE/CVF Conference on*
684 *Computer Vision and Pattern Recognition*, pp. 9579–9589, 2024.
- 685 Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman
686 Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. Bloom: A 176b-
687 parameter open-access multilingual language model. 2023.
- 688 Byung-Kwan Lee, Sangyun Chung, Chae Won Kim, Beomchan Park, and Yong Man Ro. Phantom
689 of latent for large language and vision models. *arXiv preprint arXiv:2409.14713*, 2024a.
- 690 Byung-Kwan Lee, Sangyun Chung, Chae Won Kim, Beomchan Park, and Yong Man Ro. Trol:
691 Traversal of layers for large language and vision models. *arXiv preprint arXiv:2406.12246*,
692 2024b.
- 693 Byung-Kwan Lee, Chae Won Kim, Beomchan Park, and Yong Man Ro. Meteor: Mamba-based
694 traversal of rationale for large language and vision models. *arXiv preprint arXiv:2405.15574*,
695 2024c.

- 702 Byung-Kwan Lee, Beomchan Park, Chae Won Kim, and Yong Man Ro. Collavo: Crayon large
703 language and vision model. *arXiv preprint arXiv:2402.11248*, 2024d.
- 704
- 705 Byung-Kwan Lee, Beomchan Park, Chae Won Kim, and Yong Man Ro. Moai: Mixture of all
706 intelligence for large language and vision models. *arXiv preprint arXiv:2403.07508*, 2024e.
- 707
- 708 Joonhyung Lee, Jeongin Bae, Byeongwook Kim, Se Jung Kwon, and Dongsoo Lee. To fp8 and
709 back again: Quantifying the effects of reducing precision on llm training stability. *arXiv preprint*
710 *arXiv:2405.18710*, 2024f.
- 711
- 712 Seongyun Lee, Sue Hyun Park, Seungone Kim, and Minjoon Seo. Aligning to thousands of prefer-
713 ences via system message generalization. *arXiv preprint arXiv:2405.17977*, 2024g.
- 714
- 715 Young-Jun Lee, Chae-Gyun Lim, Yunsu Choi, Ji-Hui Lm, and Ho-Jin Choi. Personachatgen: Gen-
716 erating personalized dialogues using gpt-3. In *Proceedings of the 1st Workshop on Customized*
Chat Grounding Persona and Knowledge, pp. 29–48, 2022.
- 717
- 718 Young-Jun Lee, Jonghwan Hyeon, and Ho-Jin Choi. Large language models can share images, too!
719 *arXiv preprint arXiv:2310.14804*, 2023.
- 720
- 721 Young-Jun Lee, Dokyong Lee, Junyoung Youn, Kyeongjin Oh, Byungsoo Ko, Jonghwan Hyeon,
722 and Ho-Jin Choi. Stark: Social long-term multi-modal conversation with persona commonsense
knowledge. *arXiv preprint arXiv:2407.03958*, 2024h.
- 723
- 724 Karel Lenc and Andrea Vedaldi. Understanding image representations by measuring their equiv-
725 ariance and equivalence. In *Proceedings of the IEEE conference on computer vision and pattern*
726 *recognition*, pp. 991–999, 2015.
- 727
- 728 Baiqi Li, Zhiqiu Lin, Wenxuan Peng, Jean de Dieu Nyandwi, Daniel Jiang, Zixian Ma, Simran
729 Khanuja, Ranjay Krishna, Graham Neubig, and Deva Ramanan. Naturalbench: Evaluating vision-
language models on natural adversarial samples. *arXiv preprint arXiv:2410.14669*, 2024a.
- 730
- 731 Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei
732 Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint*
733 *arXiv:2408.03326*, 2024b.
- 734
- 735 Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Bench-
736 marking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*,
2023.
- 737
- 738 Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li.
739 Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv*
740 *preprint arXiv:2407.07895*, 2024c.
- 741
- 742 Jiaang Li, Yova Kementchedjhieva, Constanza Fierro, and Anders Søgaard. Do vision and language
743 models share concepts? a vector space alignment study, 2024d. URL <https://arxiv.org/abs/2302.06555>.
- 744
- 745 Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng
746 Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models.
747 *arXiv preprint arXiv:2403.18814*, 2024e.
- 748
- 749 Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization*
branches out, pp. 74–81, 2004.
- 750
- 751 Shanchuan Lin, Anran Wang, and Xiao Yang. Sdxl-lightning: Progressive adversarial diffusion
752 distillation. *arXiv preprint arXiv:2402.13929*, 2024.
- 753
- 754 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
755 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer*
Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014,
Proceedings, Part V 13, pp. 740–755. Springer, 2014.

- 756 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction
757 tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
758 pp. 26296–26306, 2024a.
- 759 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.
760 Llava-next: Improved reasoning, ocr, and world knowledge, 2024b.
- 761 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances*
762 *in neural information processing systems*, 36, 2024c.
- 763 Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Lam Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-
764 tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks.
765 *arXiv preprint arXiv:2110.07602*, 2021.
- 766 Ziyu Liu, Tao Chu, Yuhang Zang, Xilin Wei, Xiaoyi Dong, Pan Zhang, Zijian Liang, Yuanjun Xiong,
767 Yu Qiao, Dahua Lin, et al. Mmdu: A multi-turn multi-image dialog understanding benchmark
768 and instruction-tuning dataset for lvlms. *arXiv preprint arXiv:2406.11833*, 2024d.
- 769 Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren,
770 Zhuoshu Li, Yaofeng Sun, et al. Deepseek-vl: towards real-world vision-language understanding.
771 *arXiv preprint arXiv:2403.05525*, 2024.
- 772 Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,
773 Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for
774 science question answering. In *The 36th Conference on Neural Information Processing Systems*
(*NeurIPS*), 2022a.
- 775 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-
776 Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of
777 foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- 778 Yujie Lu, Wanrong Zhu, Xin Eric Wang, Miguel Eckstein, and William Yang Wang. Imagination-
779 augmented natural language understanding. *arXiv preprint arXiv:2204.08535*, 2022b.
- 780 Xu Ma, Can Qin, Haoxuan You, Haoxi Ran, and Yun Fu. Rethinking network design and local
781 geometry in point cloud: A simple residual mlp framework. *arXiv preprint arXiv:2202.07123*,
782 2022.
- 783 Xu Ma, Yuqian Zhou, Huan Wang, Can Qin, Bin Sun, Chang Liu, and Yun Fu. Image as set of
784 points. *arXiv preprint arXiv:2303.01494*, 2023.
- 785 Subhansu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained
786 visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- 787 Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A bench-
788 mark for question answering about charts with visual and logical reasoning. *arXiv preprint*
789 *arXiv:2203.10244*, 2022.
- 790 Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter,
791 Dhruvi Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. Mm1: Methods, analysis & insights
792 from multimodal llm pre-training. *arXiv preprint arXiv:2403.09611*, 2024.
- 793 Jack Merullo, Louis Castricato, Carsten Eickhoff, and Ellie Pavlick. Linearly mapping from image
794 to text space. *arXiv preprint arXiv:2209.15162*, 2022.
- 795 Muhammad Muzammal Naseer, Kanchana Ranasinghe, Salman H Khan, Munawar Hayat, Fahad
796 Shahbaz Khan, and Ming-Hsuan Yang. Intriguing properties of vision transformers. *Advances in*
797 *Neural Information Processing Systems*, 34:23296–23308, 2021.
- 798 Yasumasa Onoe, Sunayana Rane, Zachary Berger, Yonatan Bitton, Jaemin Cho, Roopal Garg,
799 Alexander Ku, Zarana Parekh, Jordi Pont-Tuset, Garrett Tanzer, et al. Docci: Descriptions of
800 connected and contrasting images. *arXiv preprint arXiv:2404.19753*, 2024.
- 801 OpenAI. ChatGPT. <https://openai.com/blog/chatgpt/>, 2023.
- 802
- 803
- 804
- 805
- 806
- 807
- 808
- 809

- 810 OpenAI. Gpt-4v(ision) technical work and authors, 2023. [https://openai.com/](https://openai.com/contributions/gpt-4v)
811 [contributions/gpt-4v](https://openai.com/contributions/gpt-4v), Last accessed on 2024-02-13.
812
- 813 Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012*
814 *IEEE conference on computer vision and pattern recognition*, pp. 3498–3505. IEEE, 2012.
- 815 Jordi Pont-Tuset, Jasper Uijlings, Soravit Changpinyo, Radu Soricut, and Vittorio Ferrari. Con-
816 necting vision and language with localized narratives. In *Computer Vision–ECCV 2020: 16th*
817 *European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pp. 647–664.
818 Springer, 2020.
- 819 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
820 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
821 models from natural language supervision. In *International conference on machine learning*, pp.
822 8748–8763. PMLR, 2021.
- 823
- 824 Pooyan Rahmanzadehgervi, Logan Bolton, Mohammad Reza Taesiri, and Anh Totti Nguyen. Vision
825 language models are blind. *arXiv preprint arXiv:2407.06581*, 2024.
- 826
- 827 Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl:
828 Incremental classifier and representation learning. In *Proceedings of the IEEE conference on*
829 *Computer Vision and Pattern Recognition*, pp. 2001–2010, 2017.
- 830 Min Shi, Fuxiao Liu, Shihao Wang, Shijia Liao, Subhashree Radhakrishnan, De-An Huang, Hongxu
831 Yin, Karan Sapra, Yaser Yacoob, Humphrey Shi, et al. Eagle: Exploring the design space for
832 multimodal llms with mixture of encoders. *arXiv preprint arXiv:2408.15998*, 2024.
- 833 Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for
834 image captioning with reading comprehension. In *Computer Vision–ECCV 2020: 16th European*
835 *Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pp. 742–758. Springer,
836 2020.
- 837 Krishna Srinivasan, Karthik Raman, Jiecao Chen, Michael Bendersky, and Marc Najork. Wit:
838 Wikipedia-based image text dataset for multimodal multilingual machine learning. In *Proceed-*
839 *ings of the 44th international ACM SIGIR conference on research and development in information*
840 *retrieval*, pp. 2443–2449, 2021.
- 841
- 842 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya
843 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open
844 models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- 845 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide
846 shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF*
847 *Conference on Computer Vision and Pattern Recognition*, pp. 9568–9578, 2024.
- 848
- 849 Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and
850 Hervé Jégou. Training data-efficient image transformers & distillation through attention. In
851 *International conference on machine learning*, pp. 10347–10357. PMLR, 2021.
- 852
- 853 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
854 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
855 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 856
- 857 Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image
858 description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern*
859 *recognition*, pp. 4566–4575, 2015.
- 860
- 861 Gaurav Verma, Minje Choi, Kartik Sharma, Janelle Watson-Daniels, Sejoon Oh, and Srijan Kumar.
862 Cross-modal projection in multimodal llms doesn’t really project visual attributes to textual space.
863 In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*
(*Volume 2: Short Papers*), pp. 657–664, 2024.
- Kirill Vishniakov, Zhiqiang Shen, and Zhuang Liu. Convnet vs transformer, supervised vs clip:
Beyond imagenet accuracy. *arXiv preprint arXiv:2311.09215*, 2023.

- 864 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,
865 Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the
866 world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- 867 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
868 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
869 *neural information processing systems*, 35:24824–24837, 2022.
- 870 Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Chunyi Li,
871 Wenxiu Sun, Qiong Yan, Guangtao Zhai, et al. Q-bench: A benchmark for general-purpose
872 foundation models on low-level vision. *arXiv preprint arXiv:2309.14181*, 2023.
- 873 Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin,
874 Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer
875 to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024.
- 876 Hu Xu, Saining Xie, Xiaoqing Ellen Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen
877 Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Demystifying clip data. *arXiv*
878 *preprint arXiv:2309.16671*, 2023.
- 879 Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj
880 Prabhu, Yutong Dai, Michael S Ryoo, et al. xgen-mm (blip-3): A family of open large multimodal
881 models. *arXiv preprint arXiv:2408.08872*, 2024.
- 882 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
883 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*
884 *arXiv:2407.10671*, 2024.
- 885 Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang,
886 and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv*
887 *preprint arXiv:2308.02490*, 2023.
- 888 Youngjoon Yu, Sangyun Chung, Byung-Kwan Lee, and Yong Man Ro. Spark: Multi-vision sen-
889 sor perception and reasoning benchmark for large-scale vision-language models. *arXiv preprint*
890 *arXiv:2408.12114*, 2024.
- 891 Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruysen, Carlos Riquelme, Mario
892 Lucic, Josip Djolonga, Andre Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, et al. A
893 large-scale study of representation learning with the visual task adaptation benchmark. *arXiv*
894 *preprint arXiv:1910.04867*, 2019.
- 895 Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language
896 image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer*
897 *Vision*, pp. 11975–11986, 2023.
- 898 Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. Investigating
899 the catastrophic forgetting in multimodal large language model fine-tuning. In *Conference on*
900 *Parsimony and Learning*, pp. 202–227. PMLR, 2024.
- 901 Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. Large language models
902 are not robust multiple choice selectors. In *The Twelfth International Conference on Learning*
903 *Representations*, 2023a.
- 904 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
905 Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
906 Judging llm-as-a-judge with mt-bench and chatbot arena, 2023b.
- 907 Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia
908 Efrat, Ping Yu, Lili Yu, et al. Lima: Less is more for alignment. *Advances in Neural Information*
909 *Processing Systems*, 36, 2024.
- 910 Xuhui Zhou, Hao Zhu, Leena Mathur, Ruohong Zhang, Haofei Yu, Zhengyang Qi, Louis-Philippe
911 Morency, Yonatan Bisk, Daniel Fried, Graham Neubig, et al. Sotopia: Interactive evaluation for
912 social intelligence in language agents. *arXiv preprint arXiv:2310.11667*, 2023a.

918 Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit
 919 Bansal, and Huaxiu Yao. Analyzing and mitigating object hallucination in large vision-language
 920 models. *arXiv preprint arXiv:2310.00754*, 2023b.

923 A LIMITATIONS AND DISCUSSIONS

924
 925 In this work, we observe intriguing findings regarding LLVMs under various experimental settings.
 926 To provide a clear and well-defined scope for our conclusions, we further discuss the limitations of
 927 the experimental setup for our findings (or claims), explore the most plausible application directions
 928 based on our findings, and offer meaningful insights for future research directions for each finding.

929 A.1 LIMITATIONS

930
 931 Overall, our experiments have several limitations regarding model- and dataset-side generalizability,
 932 which are important for a more rigorous analysis. For instance, we primarily evaluate LLVMs on
 933 VQA-style tasks, including free-form and multiple-choice question types, and focus exclusively on
 934 the LLaVA family. To improve the generalizability of our findings, future work should explore
 935 experiments on other LLVMs, such as Qwen2-VL Wang et al. (2024), and extend evaluations to
 936 additional datasets (e.g., image captioning datasets). Furthermore, demonstrating the impact of
 937 model scaling would provide stronger support for our conclusions. Below, we present the specific
 938 limitations for each section.

939
 940 **Limitations: Section 3.3.** In Figure 1, obtaining the results required running computations for the
 941 full number of visual patch tokens, which is highly resource-intensive. This is especially challenging
 942 given the large number of visual patch tokens required by recent LLVMs—for example, 576 for
 943 LLaVA-1.5 and more than 1,000 for Qwen2-VL Wang et al. (2024).

944
 945 **Limitations: Section 3.4.** Synthesized images were generated using LLaVA-OneVision-7B Li
 946 et al. (2024b) with the prompt template: “Please generate a caption of this image.” and
 947 SDXL-Lightning Lin et al. (2024). To improve robustness, future experiments should explore
 948 captions with varying levels of detail, from concise to highly detailed, by using alternative prompt
 949 templates, specialized captioning models (e.g., ShareCaptioner⁴ Chen et al. (2023a)), or more
 950 advanced text-to-image generation models that outperform SDXL-Lightning. Incorporating these
 951 variations would enhance the reliability of our conclusions.

952
 953 **Limitations: Section 3.5.** During patch-dropping, we employed the dino-small Caron et al.
 954 (2021) model for both SaliEnt PatchDrop and Non-SaliEnt PatchDrop. The impact of patch
 955 dropping is likely to vary depending on the size and type of self-supervised vision model used (e.g.,
 956 large-scale DINO), potentially leading to differing patterns of performance degradation.

957
 958 **Limitations: Section 3.6.** While we evaluated visual perception capabilities across various image
 959 datasets, many domain-specific image datasets exist in the real world. To draw more generaliz-
 960 able conclusions, it would be beneficial to evaluate additional datasets, such as the VTAB bench-
 961 mark Zhai et al. (2019). Additionally, we investigated *catastrophic forgetting* by following existing
 962 experimental setups from the prior study Zhai et al. (2024). However, comparing LLVMs with
 963 contrastive approaches (e.g., CLIP) may be unfair due to multiple factors influencing LLVM per-
 964 formance, such as prompt variations and methods for calculating accuracy from the generated text. To
 965 enable a more rigorous analysis, future work should explore different prompt methods and fine-tune
 966 LLVMs on zero-shot image classification datasets (e.g., CIFAR-100) to assess whether perception
 967 capabilities improve. Regarding the LLM-dominance problem during visual instruction tuning, con-
 968 firming this phenomenon is challenging. To test it effectively, LLVMs should be trained with identi-
 969 cal datasets but varying LLM sizes and vision encoder scales. Alternatively, other types of LLVMs
 970 that incorporate external computer-vision models (e.g., segmentation models) such as MoAI Lee
 971 et al. (2024e) could be evaluated. Using visually enhanced LLVMs would strengthen this argument.
 In addition, for Figure 7, evaluating cross-modal alignment on a broader variety of datasets, such

⁴<https://huggingface.co/Lin-Chen/ShareCaptioner>

972 as CC12M Changpinyo et al. (2021), WIT Srinivasan et al. (2021), and RedCaps12M Desai et al.
 973 (2021), would provide a better understanding of the findings. Expanding this evaluation to various
 974 LLMs, such as LLaVA-OneVision and Qwen2-VL, would also yield more comprehensive insights.
 975

976 **Limitations: Section 3.7.** In Figures 8 and 9, obtaining the importance scores is computationally
 977 expensive. For a single run, we calculate the importance scores for each group-wise position (e.g.,
 978 36 positions for LLaVA-1.5), and we repeat the experiment K times (with $K = 10$). This results in
 979 a total of 360 experiments per benchmark. Similarly, the computation for layer importance is also
 980 resource-intensive.

981 A.2 DISCUSSIONS

982 Here, we present several discussions based on our findings.

983
 984 **Findings: Permutation Invariance.** We suggest that future work focuses on two key directions.
 985 First, it is essential to develop more challenging benchmarks that better explore LLMs’ capabilities.
 986 Such benchmarks should prioritize free-form question types and avoid including “blind” samples Fu
 987 et al. (2024); Li et al. (2024a) that models can solve using commonsense reasoning without actually
 988 perceiving the image. Building multi-turn interactive conversation benchmarks, like MMDU Liu
 989 et al. (2024d), could be particularly useful in this context. Second, since LLMs generally exhibit
 990 permutation invariance, visual patch tokens can be treated as independent elements, allowing images
 991 to be represented as unordered sets of points. Applying paradigms like “Context Clusters,” Ma
 992 et al. (2023) which rely on clustering algorithms rather than convolutions or attention mechanisms,
 993 could improve interpretability and training efficiency. Furthermore, this approach could facilitate
 994 generalization to other data domains, such as point clouds Ma et al. (2022), RGB-D data, or sensory
 995 images Yu et al. (2024), broadening the applicability of LLMs.
 996

997
 998 **Findings: Sensitivity to Spatial Structures.** One future direction is to develop more robust
 999 LLMs that can handle spatial disruptions. Real-world images often lack perfect clarity—details
 1000 may be missing, images may be flipped, or other disruptions may occur. To address this, we pro-
 1001 pose incorporating randomly shuffled images into the training process. By framing this as a jigsaw
 1002 puzzle Chen et al. (2023b) task, models can be trained to reconstruct the original positions of im-
 1003 age patches. This approach could enhance their robustness to spatial variations, making them more
 1004 applicable to real-world scenarios.

1005 **Findings: Catastrophic Forgetting.** Balancing perception and cognitive reasoning capabilities is
 1006 critical. The “catastrophic forgetting” problem Kirkpatrick et al. (2017) has been a long-standing
 1007 issue in machine learning. A standard approach is to train models on mixed datasets Ke et al. (2020);
 1008 Gururangan et al. (2020) with a carefully designed balance (a “golden ratio”) between perception-
 1009 and reasoning-related data. Continuously training LLMs on perception-focused datasets following
 1010 rehearsal methods Rebuffi et al. (2017) can minimize catastrophic forgetting by retaining knowledge
 1011 of prior tasks while learning new ones. Knowledge distillation Jin et al. (2021) from large-scale
 1012 LLMs (e.g., 72B parameters) to smaller-scale models (e.g., 7B parameters) could help preserve
 1013 perception capabilities while maintaining reasoning strength. Alternatively, fine-tuning adapters
 1014 (e.g., p-tuning Liu et al. (2021), LoRA Hu et al. (2021), Q-LoRA Dettmers et al. (2024)) on task-
 1015 specific datasets offers a lightweight solution to improve performance on new tasks without sacri-
 1016 ficing existing capabilities.

1017 **Findings: Cross-modal Alignment in the Platonic Representation Hypothesis.** Maintaining
 1018 the original cross-modal alignment is critically important. Continual learning methods (presented
 1019 above) could be applied to mitigate the loss of alignment during visual instruction tuning. Enhanc-
 1020 ing the visual perception capability of the projector during training could also help. For instance,
 1021 employing models such as HoneyBee Cha et al. (2024), which incorporate convolution layer-based
 1022 projectors, could improve localized understanding. Convolution layers are well-known for their
 1023 strong inductive bias toward localized feature extraction, making them better suited for capturing
 1024 fine-grained details in images. Even with the inclusion of complex instruction datasets (e.g., charts,
 1025 math), a carefully designed projector that excels at extracting detailed and localized information
 from images could naturally improve both perception and reasoning capabilities. We hypothesize

that enhancing localized perception would inherently lead to improvements in reasoning, aligning the two capabilities more effectively.

Findings: Importance of Central Visual Tokens. Based on our observations, reducing redundant visual tokens in the projector could enhance training and inference efficiency, aligning with findings from prior studies Alayrac et al. (2022); Cha et al. (2024); Xue et al. (2024). Typically, the large number of visual tokens poses a computational burden. This is particularly relevant for real-world scenarios where interleaved format-style conversations Li et al. (2024c); Lee et al. (2024h) are predominant. High visual token counts can make it challenging to train more effective LLVMs for such interleaved conversational formats. Our findings provide a practical direction for reducing visual token counts while maintaining performance. By doing so, we can enable the training of interleaved-format LLVM models more efficiently, similar to approaches highlighted in previous research Xue et al. (2024).

Findings: Importance of Lower Layer. Based on our observations, we emphasize the importance of the traversing layers (TroL) approach Lee et al. (2024b), in improving generalization. In this approach, models are trained to revisit and leverage layer-specific information during the training process. The paper demonstrates that lower layers are more actively engaged, which aligns with our findings. These results suggest that the lower layers of LLVMs play a critical role in establishing a foundational understanding of the world. To enhance this capability, increasing the signal for world understanding in the lower layers during training could be a promising direction. One potential method is injecting noise information into the lower layers during training, as suggested in a prior study Jain et al. (2023). This technique could improve the robustness of LLVMs, further solidifying their foundational perception and reasoning capabilities.

Findings: Relative Importance of Modalities. While the textual modality appears more influential in higher layers, improving the visual perception capability in lower layers is crucial. This is because LLVMs rely heavily on understanding the given image during the initial processing stages. As suggested in prior works Cha et al. (2024); McKinzie et al. (2024), using a larger number of visual tokens, adopting high-resolution image processing Li et al. (2024c), or employing dynamic image processing methods Wang et al. (2024); Li et al. (2024b) is essential for enhancing performance. Furthermore, strengthening the projector’s capability for localized visual understanding Cha et al. (2024) could be beneficial. For instance, after the initial image-caption alignment step (commonly the first step in LLVM training), an additional training phase called “empowering localized understanding” could be introduced before visual instruction tuning. This phase would involve adding an extra layer, referred to as the “AL” (Augmented Layer), on top of the simple linear layer. The AL would be trained using a masked autoencoder (MAE) approach He et al. (2022), where the model learns to predict masked image patches. This process would enhance localized visual understanding, ultimately improving the balance between visual and textual modalities and boosting overall model performance.

B ADDITIONAL EXPLANATION OF PLATONIC REPRESENTATION HYPOTHESIS

In Section 3.6, we investigate how effectively a trained projector preserves cross-modal alignment, drawing on the *Platonic Representation Hypothesis* Huh et al. (2024). In this section, we provide a detailed explanation of (1) the definition of the Platonic Representation Hypothesis, (2) the alignment metric, and (3) the methodology used to measure alignment in our experiment.

B.1 DEFINITION OF THE PLATONIC REPRESENTATION HYPOTHESIS

Traditionally, different types of AI models represent the world in fundamentally different ways. For instance, when presented with the same reality (e.g., an image, as illustrated in Figure 10), self-supervised vision models might focus on shapes, colors, and optical effects — features critical to visual understanding — while LLMs might emphasize semantic meanings and syntactic structures. Recently, researchers have developed LLVMs by jointly training vision models and LLMs, encouraging them to interpret and represent the world in a more unified manner. The Platonic Representation Hypothesis posits that neural networks, trained with distinct objectives on different data

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

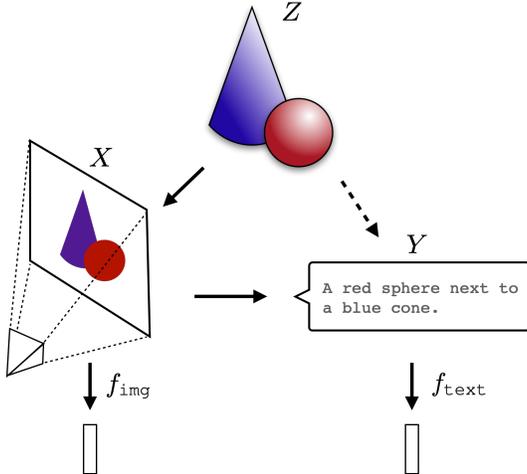


Figure 10: Images (X) and text (Y) are projections of a common underlying reality (Z). We conjecture that representation learning algorithms will converge on a shared representation of Z , and scaling model size, as well as data and task diversity, drives this convergence. For clarity, this figure and its caption have been taken exactly as they appear in the original paper Huh et al. (2024).

modalities, converge toward a shared statistical model of reality in their representation spaces. In the original paper introducing this hypothesis, the authors demonstrated a strong level of alignment between the representations of models trained on disparate modalities (e.g., Figure 3 in the original paper). Based on these findings, we argue that the alignment between models trained on different modalities should not only be preserved but potentially strengthened.

B.2 ALIGNMENT MEASUREMENT.

To evaluate the alignment between representations from two models, we employ the **Mutual k -Nearest Neighbor (MNN) Metric**. This metric focuses on local similarity by computing the intersection of the k -nearest neighbor sets for each sample from the two models’ representation spaces. The alignment is then measured based on the size of these intersections, as detailed below.

Mutual k -Nearest Neighbor Metric. Let f and g denote the representation functions of two models, and let \mathcal{X} represent the data distribution (e.g., an image-caption dataset).

1. The representations for a mini-batch of samples $\{x_i, y_i\}_{i=1}^b$ are defined as:

$$\phi_i = f(x_i), \quad \psi_i = g(y_i), \quad i = 1, \dots, b,$$

where $\Phi = \{\phi_1, \dots, \phi_b\}$ and $\Psi = \{\psi_1, \dots, \psi_b\}$ represent the feature sets produced by models f and g , respectively.

2. For each feature ϕ_i and ψ_i , the k -nearest neighbor sets are computed as:

$$\mathcal{S}(\phi_i) = \{k \text{ nearest neighbors of } \phi_i\}, \quad \mathcal{S}(\psi_i) = \{k \text{ nearest neighbors of } \psi_i\}.$$

3. The alignment for a given pair of features (ϕ_i, ψ_i) is defined as the normalized size of the intersection of their k -nearest neighbor sets:

$$m_{\text{NN}}(\phi_i, \psi_i) = \frac{1}{k} |\mathcal{S}(\phi_i) \cap \mathcal{S}(\psi_i)|,$$

where $|\cdot|$ represents the size of the intersection.

4. The overall alignment for the mini-batch is computed as the average alignment across all samples:

$$M_{\text{NN}} = \frac{1}{b} \sum_{i=1}^b m_{\text{NN}}(\phi_i, \psi_i).$$

B.3 HOW TO MEASURE IN OUR EXPERIMENT

To assess the alignment between a suite of large language models (LLMs) and vision models, we utilize the image-caption pair dataset DOCCI Onoe et al. (2024). Specifically, in DOCCI, the dataset consists of image-caption pairs

$$D = \{(x_i, y_i)\}_{i=1}^{|D|},$$

where x_i denotes the image and y_i denotes the corresponding caption text.

For our experiment, we prepare three models: an LLM (f_L), a vision encoder from a vision-language model without visual instruction tuning (f_V), and a vision encoder with a projector, representing a vision-language model with visual instruction tuning (f_{VP}). The vision encoder in f_{VP} is kept identical to f_V . For example, CLIP-L/336px is used as the vision encoder for both f_V and f_{VP} when paired with LLaVA-1.5.

In our experiment, we explore the degree of alignment lost after visual instruction tuning, guided by the Platonic representation hypothesis. We assume that in a successful LLVM, the projector should effectively represent the visual world and enable the LLM to understand and interpret the given image accurately. We calculate two alignment scores: one between f_L and f_V , and another between f_L and f_{VP} . The discrepancy between these scores reflects the extent to which alignment performance deteriorates.

To compute the alignment scores, we follow these steps:

1. Extract features from f_L by providing the input text y_i . We then apply average pooling to all the extracted hidden states.
2. Extract features from f_V by providing the image x_i , using only the feature corresponding to the [CLS] token.
3. Extract features from f_{VP} by providing the image x_i , applying average pooling to all visual patch tokens (e.g., 576 tokens in LLaVA-1.5) produced by the projector.

Finally, we calculate the alignment scores using these extracted features via the mutual nearest-neighbor alignment metric.

B.4 MOTIVATION BEHIND SELECTING THE DOCCI DATASET

We posit that the ability to perceive and reason based on complex images (e.g., charts, mathematical representations, code snippets, and diagrams) is crucial for creating a helpful assistant. However, we believe that an LLVM must first excel at understanding more natural scenes to become a broadly applicable personal AI assistant, such as one integrated into smart glasses (e.g., Meta AI’s glasses⁵) or real-time cameras (e.g., Project Astra⁶). To achieve effective alignment between the language and vision modalities, we require paired datasets where the captions provide detailed descriptions of the corresponding images. These descriptions must include essential visual features such as attributes, spatial relationships, object counts, objects, text rendering, viewpoints, optical effects, and world knowledge. Based on this criterion, we sought an image-caption pair dataset emphasizing (1) natural scenes and (2) highly descriptive captions. The DOCCI dataset meets these requirements effectively. Of course, other datasets could also be considered as candidates, such as Localized Narratives Pont-Tuset et al. (2020), CC12M Changpinyo et al. (2021), COCO-Caption Lin et al. (2014), WIT Srinivasan et al. (2021), or RedCaps12M Desai et al. (2021). In future work, we plan to conduct additional experiments to enhance the generalizability of our observations.

C ADDITIONAL EXPLANATION OF IMPORTANCE SCORE

In Section 3.7, we investigate the model’s behavior to assess the importance of either a specific layer or a visual token when performing downstream tasks. We hypothesize that introducing arbitrary noise to a specific component — either a layer block or a visual token — will significantly drop the model’s performance if that component is crucial to the reasoning process. To quantify this,

⁵<https://www.meta.com/smart-glasses/>

⁶<https://www.youtube.com/watch?v=nXVvRhiGjI>

we define an *importance score* (\mathcal{I}), inspired by the concept of “sharpness of minima.” This section provides a detailed explanation of how the importance score is computed.

How is Arbitrary Noise Introduced into Target Layers or Visual Tokens? Based on Equation (2), we prepare the constraint candidate set \mathcal{C}_t , defined as a squared boundary:

$$-\epsilon + |x_t| \leq z_t \leq \epsilon + |x_t|, \quad (4)$$

where $\epsilon \sim \text{Uniform}(-1, 1)$. At each iteration, we randomly sample a noise vector z_t and apply it to the target component. Below, we detail how this is done for visual tokens, layers, and modalities.

- Visual Token Importance:** When evaluating the importance of a visual patch token (Figure 8), the noise vector is injected into the group-wise visual patch token embeddings at the target position. For instance, Figure 8 illustrates 36 positions. To measure the importance of position 0, we add the noise vector to the corresponding visual patch token embeddings at position 0, while leaving all other patch token embeddings unchanged. These modified embeddings are then input into the LLM for further processing.
- Layer-Wise Importance:** To explore layer-wise importance, the noise vector is injected into the target layer before it is processed by the LLM. Specifically, the noise is applied directly to the layer’s input embeddings before passing the target layer, ensuring that the perturbation affects only the selected layer.
- Modality Importance:** To calculate the importance of the textual modality (\mathcal{I}_T), the noise vector is injected only into the positions corresponding to text inputs within the target layer, while leaving the positions associated with visual patch tokens unchanged. Conversely, for visual modality importance (\mathcal{I}_I), the noise vector is injected into the positions corresponding to visual patch tokens within the target layer. The relative importance score for each modality is then computed as $\frac{\mathcal{I}_I}{\mathcal{I}_T}$.

To enable better interpretation across layers, all importance scores (both layer-wise and modality-specific) are normalized using min-max normalization.

D ADDITIONAL EXPLANATION OF EXPERIMENTAL SETUP

In this section, we provide a more detailed explanation of the experimental setup used to obtain our findings, including the required models, preparation of corrupted images, and other specifics. All experiments were conducted using eight A100 GPUs (40GB).

Experimental Setup: Section 3.3. We prepared ViT-variant vision encoder-equipped LVLMS that incorporate visual patch tokens. The experiments focus on visual patch tokens processed after the projector. Before conducting the “permutation invariance” experiments, we first demonstrated whether each visual patch token contains localized information. For the experiment on “sensitivity to spatial structure,” shuffled images were used, as shown in Figure 2, following the methodology of a prior study Naseer et al. (2021).

Experimental Setup: Section 3.4. To generate synthesized images, we utilized an image captioner (llava-hf/llava-onevision-qwen2-7b-ov-hf) combined with a text-to-image generative model (sdxl_lightning_8step_unet.safetensors). Additionally, a prompt template was carefully designed for this purpose.

Experimental Setup: Section 3.5. We prepared occluded images using three masking methods as described in prior work Naseer et al. (2021): Random PatchDrop, Salient PatchDrop, and Non-Salient PatchDrop. To implement Salient PatchDrop and Non-Salient PatchDrop, we employed the dino-small model Caron et al. (2021). Furthermore, to evaluate the robustness of LVLMS to occlusion, we first verified whether ViT-variant encoders exhibit genuine robustness to occlusion by comparing them with CNN-based counterparts, such as ResNet.

LLVMs	MMVP	Q-Bench	MME	MMStar	MM-Vet	LLaVA ^W	MathVista	SQA ^I	ChartQA	AI2D	Avg. Δ
LLaVA-1.5	34.67	59.73	1850.07	34.20	31.50	67.50	24.70	65.59	16.92	53.34	
+ Perm.	36.00	59.60	1874.60	33.33	30.40	66.20	21.20	65.44	14.08	52.69	▼ 0.59
	(▲ 1.33)	(▼ 0.13)	(▲ 24.53)	(▼ 0.87)	(▼ 1.10)	(▼ 1.30)	(▼ 3.50)	(▼ 0.15)	(▼ 2.84)	(▼ 0.65)	
LLaVA-NeXT	36.67	63.55	1874.42	37.80	43.50	75.50	32.00	62.12	66.06	64.02	
+ Perm.	37.33	62.54	1890.19	36.87	43.40	75.80	21.70	62.12	34.55	64.02	▼ 2.71
	(▲ 0.67)	(▼ 1.00)	(▲ 15.78)	(▼ 0.93)	(▼ 0.10)	(▲ 0.30)	(▼ 10.30)	(▼ 0.00)	(▼ 31.51)	(▼ 0.00)	
LLaVA-OneVision	60.67	77.26	1982.5	59.87	57.80	87.40	61.80	94.00	93.52	81.25	
+ Perm.	59.33	76.99	1964.3	54.93	47.60	82.30	53.50	89.24	58.26	75.58	▼ 9.40
	(▼ 1.33)	(▼ 0.27)	(▼ 18.2)	(▼ 4.93)	(▼ 10.20)	(▼ 5.10)	(▼ 8.30)	(▼ 4.76)	(▼ 35.26)	(▼ 5.67)	
QwenVL-2	50.67	77.06	2356.70	55.27	62.60	94.10	59.80	0.00	94.83	80.21	
+ Perm.	48.67	77.19	2266.96	53.47	62.20	93.20	53.10	0.00	83.59	77.43	▼ 12.82
	(▼ 2.00)	(▲ 0.13)	(▼ 89.74)	(▼ 1.80)	(▼ 0.40)	(▼ 0.90)	(▼ 6.70)	(▼ 0.00)	(▼ 11.25)	(▼ 2.78)	
Fuyu-8B	30.00	40.33	0.00	19.67	16.30	0.00	0.00	0.00	15.81	0.00	
+ Perm.	28.67	38.80	0.00	18.93	10.90	0.00	0.00	0.00	7.50	0.00	▼ 6.92
	(▼ 1.33)	(▼ 1.54)	(▼ 0.00)	(▼ 0.73)	(▼ 5.40)	(▼ 44.00)	(▼ 0.00)	(▼ 0.00)	(▼ 8.31)	(▼ 0.00)	

Table 4: Results of drop ratio (Δ) when random permutation is applied. We run five experiments.

Experimental Setup: Section 3.6. We curated image classification datasets containing realistic and natural images across various domains. To explore the platonic representation hypothesis Huh et al. (2024), we first thoroughly examined its definition, as detailed in Appendix B. This process involved preparing a diverse set of LLMs, vision encoders, and vision encoders equipped with projectors in LVLMs. We also selected datasets for verifying cross-modal alignment, ensuring that they included natural and realistic images.

Experimental Setup: Section 3.7. We first clarified the definition of “importance score” and determined how to introduce noise into the visual patch tokens. This procedure is described in Appendix C.

E ADDITIONAL EXPERIMENTAL RESULTS

E.1 PERMUTATION INVARIANCE.

As shown in Table 4, we investigate the extent to which other LVLMs exhibit permutation invariance under the same experimental settings described in Table 1. Overall, the Qwen2-VL-7B Wang et al. (2024) and Fuyu-8B models Bavishi et al. (2023) demonstrate permutation invariance on average, displaying patterns similar to those observed in the LLaVA-family models. A more detailed analysis across benchmarks reveals interesting patterns. In perception-focused benchmarks, such as MMVP, Q-Bench, MME, and MMStar (the latter two being integrated capability benchmarks that include perception-related tasks), the performance drop due to permutation is negligible. However, in text-rich benchmarks like MathVista and ChartQA, the performance drops significantly. These benchmarks require an understanding of detailed numerical information and highly structured geometric graphs, where maintaining the spatial structure of visual patch tokens is critical.

Difficulty of Benchmark. Interestingly, in the SQA^I benchmark, which includes science-related datasets, and the AI2D benchmark, which consists of diagram images, the relatively small performance gap is noteworthy, even though these images are rich in detail. We speculate that this phenomenon might be influenced by the difficulty of the benchmark, particularly the “question type.” Benchmarks typically include two question formats: (1) free-form and (2) multiple-choice questions (MCQ). We hypothesize that:

1. LLMs can often solve questions using their extensive commonsense reasoning, even without image perception. Li et al. (2024a); Fu et al. (2024)
2. MCQ formats may be easier for models compared to free-form questions due to the presence of preferred answer patterns or inherent biases in selection.

To investigate further, we conduct additional experiments comparing the difficulty of MathVista, ChartQA, SQA^I, and AI2D. We randomly select 500 samples from each dataset and, for MCQ

Datasets	Question Type	Accuracy (%)	Don't Know (%)
MathVista	Free-Form	0.3	82.1
	MCQ	36.8	0
	Overall	13.6	52.2
ChartQA	Free-Form	0	90
SQA ^I	MCQ	64.2	0
AI2D	MCQ	53.2	1.6

Table 5: Accuracy results of ChatGPT on four benchmarks for two different question types.

samples, include only those with four options. We then prompt ChatGPT (i.e., gpt-3.5-turbo) to answer these questions using the following templates:

Prompt Template for MCQ

Question: {question}
 Choices:
 {choices}
 E: I don't know.

Please MUST generate only one option (A, B, C, D, E). Do not generate any explanation.
 Answer:

Prompt Template for Free-Form

Question: {question}

Please provide your answer. If it is difficult to provide an answer, respond with "I don't know."

We added the “*I don't know*” option to prevent the model from guessing randomly. Table 5 show that ChatGPT performs better on MCQ-type benchmarks compared to free-form types. Moreover, ChatGPT achieves higher accuracy on AI2D and SQA^I compared to MathVista and ChartQA. This supports the observation that LLMs exhibit less permutation invariance in these text-rich benchmarks, possibly due to the nature of the datasets and their question formats. For free-form questions, the “don't know” response rate is significantly higher, indicating that these benchmarks are more challenging. This highlights the need to minimize “blind” samples — questions solvable by LLMs without image perception — in benchmark design. Benchmarks should prioritize free-form questions to reduce potential selection bias Zheng et al. (2023a), as argued by recent studies Li et al. (2024a).

E.2 SENSITIVITY TO SPATIAL STRUCTURES

As shown in Figure 11, we randomly shuffle image patches to evaluate their impact on model performance and observe that Qwen2-VL exhibits a similar tendency to LLaVA-family models. Specifically, we found that Qwen2-VL and LLaVA-OneVision are highly sensitive to spatial structures in text-rich benchmarks (e.g., MathVista, AI2D), which contain detailed numerical information. Notably, the performance of the Qwen2-VL model dropped significantly when the grid size was 2. To understand why Qwen2-VL is particularly sensitive, we hypothesize that this behavior is linked to its use of enhanced multi-modal rotary position embeddings (M-ROPE) Wang et al. (2024). This embedding mechanism likely contributes to the performance degradation observed when image patches are shuffled. Conversely, the model is relatively insensitive to spatial structures in perception-centric benchmarks (e.g., MMVP).

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

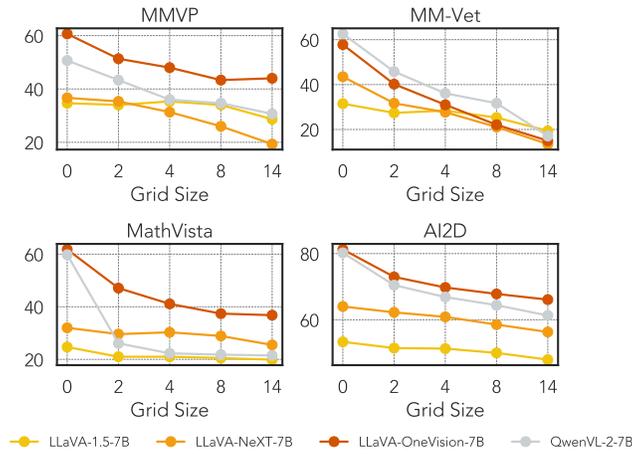


Figure 11: We present the performance across different grid sizes (2, 4, 8, 14) on the MMVP, MM-Vet, MathVista, and AI2D datasets, using four models: LLaVA-1.5, LLaVA-NeXT, LLaVA-OneVision, and Qwen2-VL.

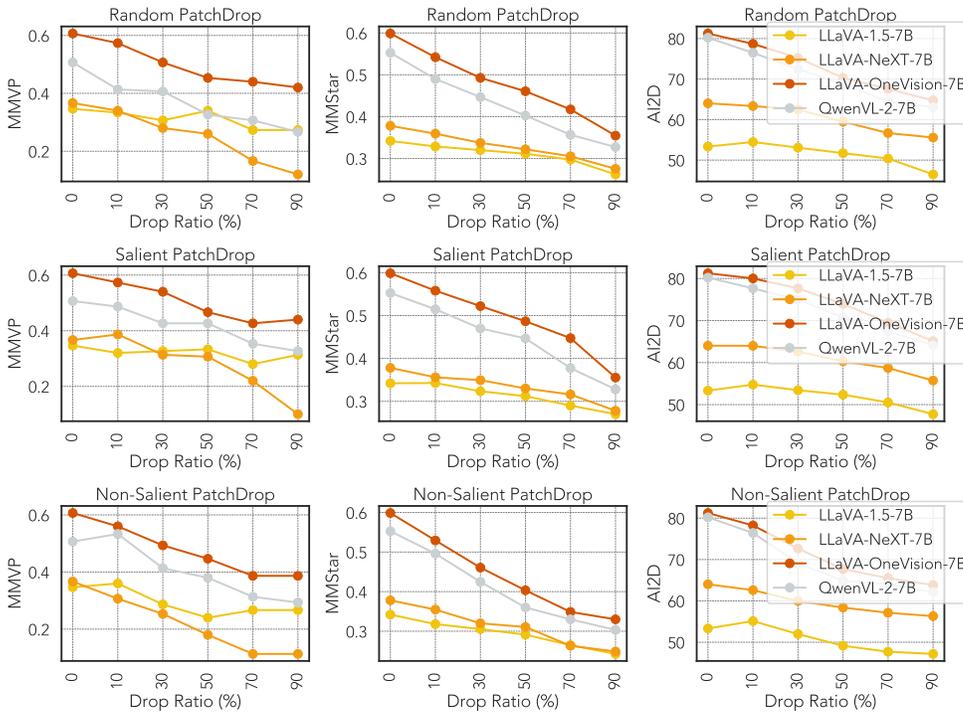


Figure 12: We present robustness performance under occlusion conditions.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

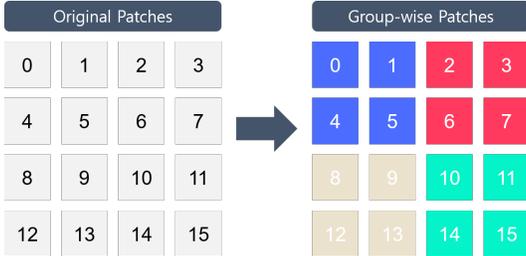


Figure 13: An illustration of group-wise patching.

E.3 OCLUSIONS

In Figure 12, we observe that the Qwen2-VL model exhibits a similar tendency to the LLaVA family models. Notably, the performance trend slope of the Qwen2-VL model closely resembles that of LLaVA-OneVision, suggesting that both models — currently high-performing LVLMs — share similar patterns. This alignment supports the generalizability of our observations. Specifically, LVLMs demonstrate relatively strong performance under occlusion. For instance, in the AI2D dataset, even when 50–70% of image patches are missing, the models can still provide correct answers to some extent. Moreover, in these scenarios, the Qwen2-VL and LLaVA-OneVision models outperform LLaVA-1.5 and LLaVA-NeXT, even when no patches are missing. These results indicate that state-of-the-art LVLMs possess strong visual understanding capabilities. This suggests that improving visual understanding during training contributes significantly to high performance and robustness against occlusion.

E.4 VARYING GRID SIZE FOR GROUPING STRATEGY

In Figure 1 and Figure 8, we group the nearest patches. For clarification, we visualize how the patches are grouped, as shown in Figure 13. Similar to the operation of a convolution layer, we group neighboring patches into a single group (indicated by the same color) and feed these groups into the model. Here, we vary the grid size, which corresponds to changing the number of elements in each group, and investigate whether the pattern observed in Figure 1 changes. We conduct additional experiments using a 3 × 3 grid of patches in Figure 14. We observe that increasing the number of grid patches leads to more precise observations. Compared to a 6 × 6 grid of patches, a 3 × 3 grid yields less precise observations. While conducting experiments on all visual patch tokens (576 for LLaVA-1.5) would provide the most precise interpretations, this approach is computationally intensive, as mentioned in Section 3.3. Therefore, we believe our chosen grid size strikes a reasonable balance for obtaining meaningful interpretations.



Figure 14: We demonstrate the extent to which group-wise visual tokens capture region-specific information (PIL) for LLaVA-1.5-7B on the MMStar (Chen et al., 2024a) and MME (Fu et al., 2023) when a 3 × 3 grid of patches. Darker regions indicate areas where the model retains more localized information for those specific groups.

E.5 DETAILED ANALYSIS OF NUMERICAL INFORMATION

As shown in the above Table 6, in overall, the Org. ratio of LLVM generating “1” in free-form question types is reduced compared to the Syn. cases. This results suggest that LVLMs can effectively interpret and understand the detailed numerical information in the given image, thereby, the phenomenon that LLVM tend to use their commonsense reasoning is reduced. However, considering the ratio of LLaVA-1.5 (44%), this ratio is not negligible. Therefore, in the future, we need to build more challenging benchmark that do not rely on the commonsense reasoning.

Additionally, we observe that most LVLMs prefer to answer “no” for yes/no question types in multiple-choice question (MCQ) formats. This suggests that, when presented with synthesized images, LVLMs struggle to solve the given questions effectively. Instead of attempting to provide an

Model	Syn.				Orig.		
	Freq. of 1	No (%)	Precision	Recall	Freq. of 1	Precision	Recall
LLaVA-1.5	81.0	64.0	49.2	36.8	44.4	59.4	43.7
LLaVA-NeXT	50.0	54.4	50.0	47.1	13.8	54.0	39.1
Meteor	9.5	82.8	55.2	18.4	7.5	78.0	36.8
LLaVA-OneVision	12.0	70.5	54.9	32.2	8.3	72.4	72.4

Table 6: Detailed analysis of the Syn. and Orig. versions of MathVista Lu et al. (2023). Precision and recall are reported for the yes/no question type.

Datasets	Prompt Template for CLIP	Prompt Template for LLM
Caltech101	a photo of a {c}.	What is the object in the image? Please answer only a single object in {class.labels}.
CIFAR-100	a photo of a {c}.	What is the object in the image? Please answer only a single object in {class.labels}.
Food101	a photo of {c}, a type of food	What is the type of food in the image? Please answer only a single type of food in {class.labels}.
Pets	a photo of a {c}, a type of pet.	What is the type of pet in the image? Please answer only a single type of pet in {class.labels}.
Country211	a photo showing the country of {c}.	What is the country in the image? Please answer only a single country in {class.labels}.
EuroSAT	a centered satellite photo of {c}.	What is the type of centered satellite in the image? Please answer only a single type of centered satellite in {class.labels}.
AirCraft	a photo of a {c}, a type of aircraft.	What is the type of aircraft in the image? Please answer only a single type of aircraft in {class.labels}.

Table 7: Prompt templates used for evaluating CLIP and LLMs on zero-shot image classification tasks. The c represents a single class label, while class.labels refers to all class labels provided by each dataset.

answer based on the limited or unclear information available in the synthesized images, LLMs tend to decline by answering “no,” leading to an increased frequency of “no” responses compared to “yes.” Furthermore, across all models, the Orig. dataset consistently yields better performance in both precision and recall. This indicates that LLMs face significant challenges in solving questions based on synthesized information. In the Syn. case, precision is consistently higher than recall, reflecting the tendency of LLMs to output “no” answers more frequently than “yes” answers. This behavior underscores the challenges LLMs face in effectively using synthesized visual information to provide accurate answers to yes/no questions.

E.6 ADDITIONAL RESULTS OF CROSS-MODAL ALIGNMENT

How to evaluate the zero-shot image classification task? To evaluate CLIP models on the zero-shot classification task, we use the prompt templates provided by CLIP-Benchmark⁷. All the prompt templates we used are presented in Table 7. For evaluating LLMs on the zero-shot image classification task, we design prompt templates inspired by those used for the CLIP model. Using these templates, the LLM predicts a single class label. Based on the LLM’s generated answer, we then use ChatGPT to verify the prediction. Specifically, we utilize the following prompt: *Please only answer the question in yes or no. Is the “Prediction” correctly predicting the right “Label”? Label: label; Prediction: outputs.* This evaluation method strictly follows the approach used in an existing study Zhai et al. (2024).

E.7 ADDITIONAL RESULTS ON IMAGE CAPTIONING TASK

The evaluation benchmarks used in our experiments primarily consist of VQA tasks, which focus on binary, multiple-choice, and free-form question types. To address whether our claim regarding “permutation invariance” generalizes to other datasets, we conduct additional experiments using image captioning tasks. These tasks inherently require “visual processing capabilities,” such as understanding attributes, viewpoints, scenes, and objects. For this investigation, we evaluate three standard datasets: COCO-Captions Lin et al. (2014) (Karpathy test set), NoCaps Agrawal et al. (2019) (validation set), and TextCaps Sidorov et al. (2020) (validation set). To generate captions, we follow the default prompting setup from LLMs-Eval⁸, which uses the prompt: *“Please carefully observe the image and come up with a caption for the image.”* We employ standard evaluation metrics — ROUGE-L Lin (2004) and CIDEr Vedantam et al. (2015) — to assess performance.

⁷<https://huggingface.co/clip-benchmark>

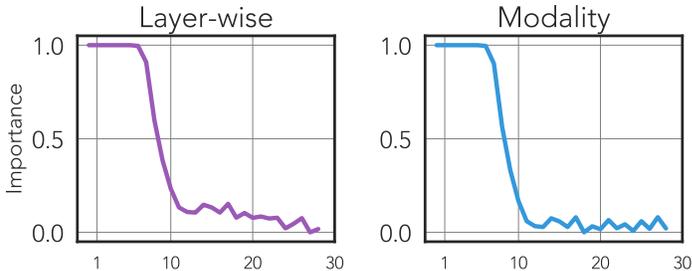
⁸<https://huggingface.co/lmms-lab>

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529

LLVMs	COCO-Captions		NoCaps		TextCaps	
	ROUGE-L	CIDEr	ROUGE-L	CIDEr	ROUGE-L	CIDEr
LLaVA-1.5	22.01	0.97	25.34	1.52	22.46	6.09
+ Perm.	22.62	1.26	26.05	2.89	22.94	7.28
Avg. Δ	\blacktriangle 0.62	\blacktriangle 0.29	\blacktriangle 0.71	\blacktriangle 1.37	\blacktriangle 0.48	\blacktriangle 1.19
LLaVA-NeXT	21.63	8.12	22.78	6.26	21.49	15.94
+ Perm.	21.86	7.64	22.68	5.81	20.19	12.30
Avg. Δ	\blacktriangle 0.24	\blacktriangledown 0.48	\blacktriangledown 0.10	\blacktriangledown 0.44	\blacktriangledown 1.29	\blacktriangledown 3.65
LLaVA-OneVision	57.23	116.25	56.09	86.60	44.58	72.69
+ Perm.	56.70	116.17	56.36	85.94	44.19	68.18
Avg. Δ	\blacktriangledown 0.53	\blacktriangledown 0.08	\blacktriangle 0.26	\blacktriangledown 0.66	\blacktriangledown 0.39	\blacktriangledown 4.52
Qwen2-VL	39.98	44.61	44.01	39.37	35.80	46.86
+ Perm.	37.19	39.29	42.70	38.35	35.31	44.64
Avg. Δ	\blacktriangledown 2.79	\blacktriangledown 5.33	\blacktriangledown 1.31	\blacktriangledown 1.02	\blacktriangledown 0.49	\blacktriangledown 2.22

1530 Table 8: Results of drop ratio (Δ) when random permutation is applied. We run five experiments.

1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541



1542 Figure 15: We present the results of (left) layer-wise importance and (right) modality importance
1543 within the layers on MME Fu et al. (2023) dataset.
1544

1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555

As shown in Table 8, we observe similar trends across image captioning datasets: **most LLMs exhibit permutation invariance**. Interestingly, on the TextCaps dataset, the performance drop is more pronounced compared to other datasets, suggesting relatively greater permutation variance. TextCaps contains more complex images (e.g., those with detailed numerical information) compared to the other datasets, which may explain this phenomenon. When comparing these findings to those in Table 1, we note that in perception-related tasks (e.g., involving natural scenes), LLMs generally exhibit permutation invariance. However, in reasoning-related tasks (e.g., MathVista) involving images with complex structures (e.g., charts or diagrams), LLMs demonstrate greater permutation variance. This suggests that maintaining the geometric or positional structure of plots and charts is crucial.

1556
1557
1558

E.8 ADDITIONAL RESULTS ON LAYER & MODALITY IMPORTANCE

1559
1560
1561
1562
1563
1564
1565

Figure 15 (left) shows that the lower layers (< 10) play a crucial role in handling integrated capabilities. Meanwhile, Figure 15 (right) demonstrates that in the lower layers (< 12), the image modality is more important than the text modality. Overall, the tendencies observed on the MME dataset are similar to those on the MMStar dataset, as shown in Figure 9. However, a key difference lies in the layer index at which the modality importance shifts; for the MME dataset, this transition occurs at a higher layer index. Based on these results, we hypothesize that LLMs allocate more effort to understanding the given images on the MME dataset compared to the MMStar dataset. One of the possible reasons is that the images in the MME dataset are more challenging for the model

to comprehend, but we can not guarantee this reason is correct, therefore, Further investigation is required to validate this assumption in future studies.

F ADDITIONAL RELATED WORKS

Model-Stitching. The model-stitching (Lenc & Vedaldi, 2015; Bansal et al., 2021) is a technique first introduced to study the internal representations of neural networks by measuring the representational similarity between two given models. Consider two models defined as $f = f^m \circ \dots \circ f^1$ and $g = g^n \circ \dots \circ g^1$. Specifically, the *stitched* model is formalized as $\mathcal{F} = g^n \circ \dots \circ g^{k+1} \circ s \circ f^k \circ \dots \circ f^1$, where s is a simple stitching layer (e.g., a linear layer or a 1×1 convolution). Therefore, even if the two models f and g differ in training methodology (e.g., supervised vs. self-supervised) or modalities (e.g., text vs. image), if \mathcal{F} exhibits good performance, then f and g have strongly correlated and compatible internal representations at layer k , apart from the stitching layer s . Merullo et al. (2022) have the similar concept of *model-stitching* to verify their strong hypothesis that the conceptual representations from a frozen LLM and a visual encoder are sufficiently similar such that a simple linear mapping layer can align them.

G ADDITIONAL EXAMPLES OF SYNTHESIZED IMAGES

We provide additional examples of synthesized images in Figure 16.

H ADDITIONAL EXAMPLES OF SHUFFLED IMAGES

We provide additional examples of shuffled images in Figure 17.

I ADDITIONAL EXAMPLES OF OCCLUDED IMAGES

We provide additional examples of occluded images in Figure 18.

J DESCRIPTION OF EVALUATION BENCHMARKS

- **MM-Vet** (Yu et al., 2023) dataset is a benchmark designed to evaluate large vision-language models (LVLMs) across six core vision-language (VL) capabilities: recognition, knowledge, optical character recognition (OCR), spatial awareness, language generation, and mathematical reasoning. The dataset includes open-ended, real-world questions based on image-text pairs, requiring models to integrate multiple capabilities to solve complex tasks. MM-Vet benchmark consists of 200 images paired with 218 open-ended questions.
- **Q-Bench** (Wu et al., 2023) evaluates the capabilities of large vision-language models in three main areas related to low-level vision tasks. These tasks focus on evaluating how well LVLMs can perform basic low-level perception tasks that are traditionally associated with human visual perception. In the Q-Bench dataset, the questions are of three types: Yes-or-No, What, and How.
 - **Low-Level Visual Perception:** Assesses how accurately LVLMs can answer questions about low-level image attributes (e.g., clarity, color, distortion). LLVisionQA dataset includes 2,990 images, each with a corresponding question about low-level features.
 - **Low-Level Visual Description:** Evaluates the ability of LVLMs to describe images. LLDescribe dataset has 499 images with expert-labeled descriptions averaging 58 words each. LVLMs are compared against these to assess completeness, preciseness, and relevance.
 - **Visual Quality Assessment:** Evaluates LVLMs’ ability to predict quantifiable quality scores for images by assessing how well they align with human-rated mean opinion scores (MOS) on low-level visual appearances, using 81,284 samples.

1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

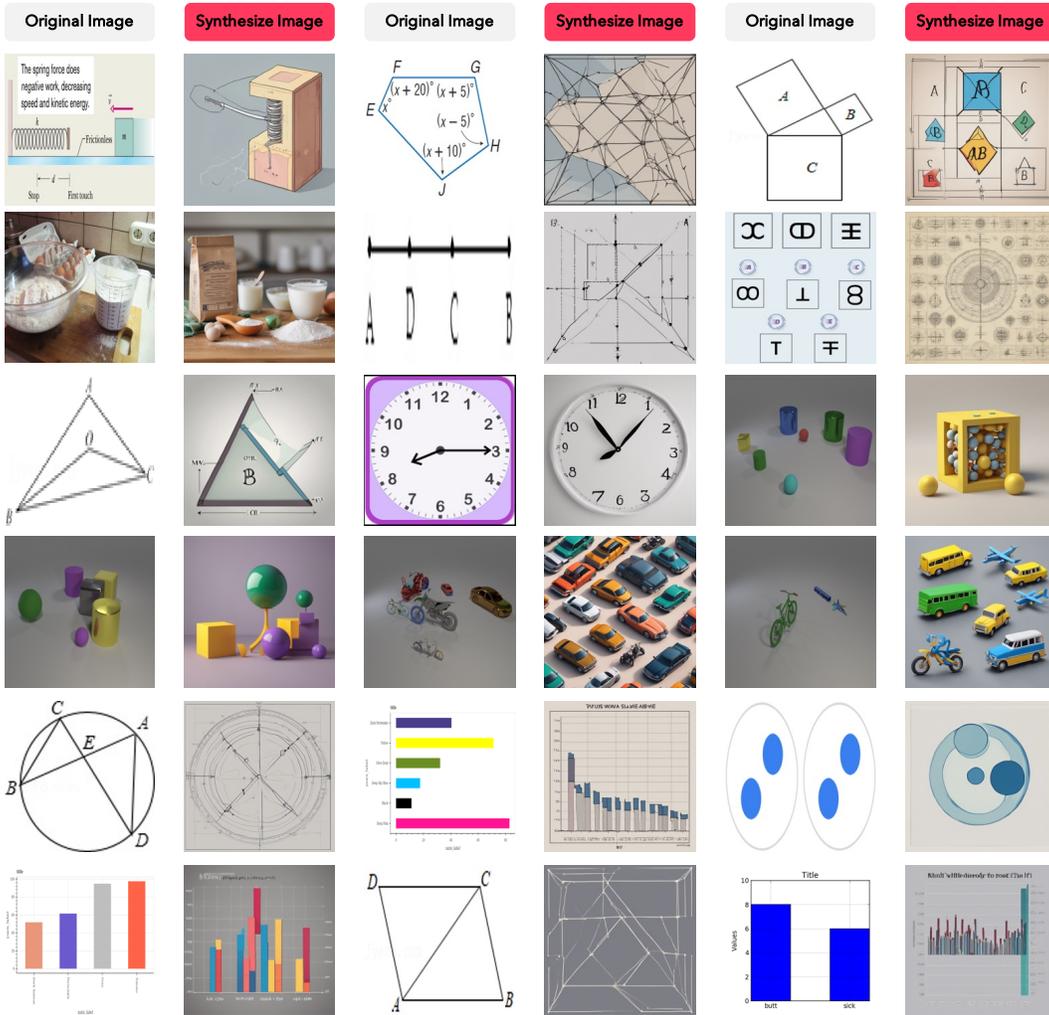


Figure 16: Examples of synthesized images from MathVista Lu et al. (2023).

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727



Figure 17: Examples of synthesized images from MM-Vet Yu et al. (2023).

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

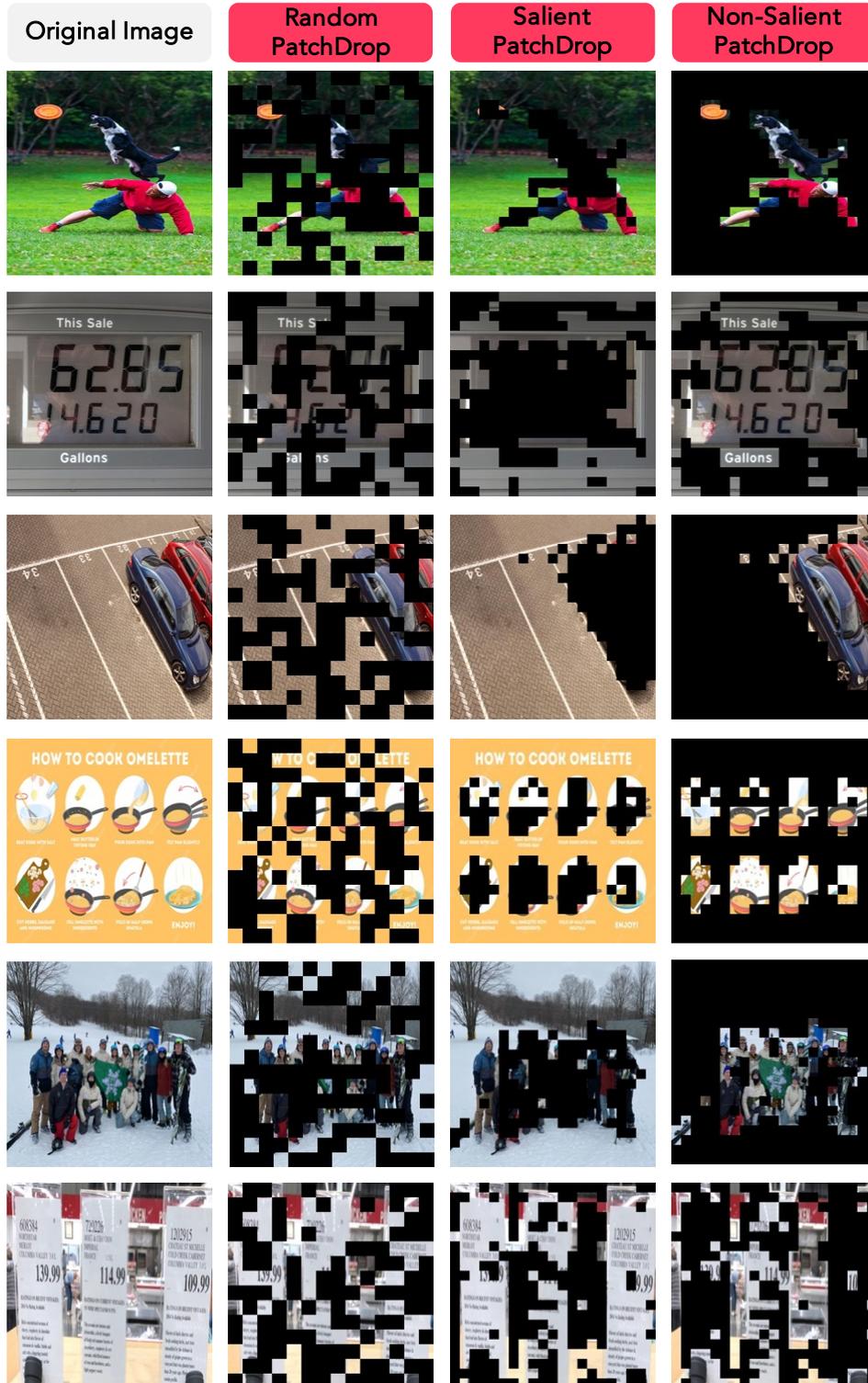


Figure 18: Examples of occluded images from MME Fu et al. (2023).

- 1782 • **SQA-IMG** (Lu et al., 2022a) is a portion of the Science Question Answering (SQA) dataset
1783 that contains questions from a wide range of scientific domains, each paired with corre-
1784 sponding image contexts. The dataset includes 10,332 examples of multimodal multiple-
1785 choice questions, along with lectures and explanations that detail the reasoning behind the
1786 correct answers.
- 1787 • **ChartQA** (Masry et al., 2022) dataset is a benchmark designed to test AI models on their
1788 ability to perform question-answering tasks over various types of charts. It focuses specifi-
1789 cally on questions requiring complex reasoning, such as visual and logical interpretation,
1790 going beyond simpler template-based datasets. ChartQA includes 9,608 human-authored
1791 open-ended questions as well as 23,111 questions that are automatically generated from
1792 chart summaries.
- 1793 • **SEED-IMG** (Li et al., 2023), a subset of SEED-Bench, focuses on evaluating spatial com-
1794 prehension of images by testing models on various dimensions like scene understanding,
1795 object identification, and spatial relationships. In terms of scale, the dataset includes 19,000
1796 multiple-choice questions that evaluate both image and video comprehension, covering 12
1797 evaluation dimensions such as scene understanding, instance identity, spatial relations, and
1798 action recognition.
- 1799 • **MME** (Fu et al., 2023) evaluates both perception and cognition abilities of LVLMS. It
1800 features 14 subtasks, including recognition tasks (such as object existence, count, position,
1801 color) and reasoning tasks (such as commonsense reasoning, numerical calculation, text
1802 translation, and code reasoning). MME uses manually created instruction-answer pairs,
1803 ensuring no overlap with public datasets. MME uses "yes/no" responses for quantitative
1804 evaluations.
- 1805 • **MathVista** (Lu et al., 2023) is a benchmark designed to evaluate the mathematical rea-
1806 soning capabilities of foundation models in visual contexts. It integrates challenges from
1807 diverse mathematical and visual tasks, with a focus on fine-grained, deep visual under-
1808 standing and compositional reasoning. MathVista consists of 6,141 examples including
1809 3,392 multiple-choice questions and 2,749 free-form questions derived from 28 existing
1810 multimodal datasets and 3 newly created datasets: IQTest, FunctionQA, and PaperQA.
- 1811 • **LLaVA-W** (Liu et al., 2024c) is a challenging evaluation benchmark created to assess the
1812 generalization and instruction-following capabilities LVLMS in complex, real-world sit-
1813 uations. It consists of 24 images and 60 questions, including diverse scenes like indoor
1814 environments, outdoor settings, memes, paintings, and sketches. Each image is associated
1815 with a highly detailed and manually curated description, and the questions focus on extract-
1816 ing intricate details and reasoning about the visual content. LLaVA-W involves a variety of
1817 tasks, including detailed descriptions, conversational answers, and complex reasoning.
- 1818 • **MMStar** (Chen et al., 2024a) is a vision-dependent multimodal benchmark designed to
1819 evaluate the multimodal capabilities of LVLMS. It addresses two main issues identified
1820 in previous benchmarks: the reliance on textual information without visual input and data
1821 leakage during training. MMStar is composed of 1,500 samples carefully selected to en-
1822 sure that visual content is necessary for solving each problem. MMStar evaluates six core
1823 capabilities across 18 detailed axes, which include tasks like image perception and logical
1824 reasoning. MMStar uses multiple-choice as the primary answer type.
- 1825 • **MMVP** (Tong et al., 2024) evaluates the visual grounding capabilities of large vision-
1826 language models by identifying scenarios where they fail to distinguish simple visual pat-
1827 terns in images. These patterns include aspects like orientation, counting, viewpoint, and
1828 relational context. The benchmark is constructed using 150 pairs of images, resulting in
1829 300 multiple-choice questions.

1830 K DESCRIPTION OF EVALUATION LVLMS

- 1832 • **LLaVA-1.5** (Liu et al., 2024a) incorporates academic task-oriented datasets to enhance
1833 performance in VQA tasks and features an MLP vision-language connector, which im-
1834 proves upon the original linear layer utilized in LLaVA (Liu et al., 2024c). It uses CLIP
1835 ViT-L/14 (Radford et al., 2021) with a 336px resolution as its vision encoder, resulting in
a total of $(336/14)^2 = 576$ visual tokens. LLaVA-1.5 is built on Vicuna with either 7B or

- 1836 13B parameters. The training dataset includes 558K samples for pre-training and 665K for
1837 fine-tuning, totaling 1.2M image-text pairs from publicly available datasets
- 1838 • **LLaVA-NeXT** (Liu et al., 2024b) (also known as LLaVA-1.6) enhances visual reasoning,
1839 OCR, and world knowledge, offering four times higher image resolution (up to 1344x336)
1840 and improved performance in visual conversations. Its architecture includes a CLIP ViT-
1841 L/14 as a vision encoder, paired with Vicuna models ranging from 7B to 34B as a back-
1842 bone language model. It utilizes 1.3M visual instruction tuning data samples for training,
1843 maintaining efficiency with approximately one day of training on 32 A100 GPUs. The archi-
1844 tecture’s high resolution and dynamic grid scheme improve detailed image processing
1845 capabilities.
 - 1846 • **LLaVA-OneVision** (Li et al., 2024c) is a LVLM designed for task transfer across single-
1847 image, multi-image, and video scenarios, with strong capabilities in video understand-
1848 ing through image-to-video task transfer. Its architecture consists of a Qwen2 language
1849 model (Yang et al., 2024) with 8B to 72B parameters, and the SigLIP vision encoder (Zhai
1850 et al., 2023), which processes images at a base resolution of 384x384, producing 729 visual
1851 tokens. The model employs a 2-layer MLP projector. The training utilized 3.2M single-
1852 image data samples and 1.6M multi-modal data samples, focusing on high-quality visual
1853 instruction tuning data to enhance its multimodal capabilities.
 - 1854 • **Meteor** (Lee et al., 2024c) is a large vision-language model that uniquely embeds multi-
1855 faceted rationales using a Mamba-based architecture (Gu & Dao, 2023), enabling efficient
1856 processing of lengthy rationales to enhance its vision-language understanding. This ap-
1857 proach allows Meteor to achieve superior performance without scaling up model size or
1858 employing additional vision encoders. Its architecture includes a CLIP-L/14 vision en-
1859 coder with an image resolution of 490x490, comprising 428M parameters, and InternLM2-
1860 7B (Cai et al., 2024) as a foundational LLM. Meteor was trained on 2.1M question-answer
1861 pairs, with 1.1M curated triples.
 - 1862 • **TroL** (Lee et al., 2024b) uses a unique characteristic called layer traversing, which reuses
1863 layers in a token-wise manner, allowing it to simulate retracing the answering process with-
1864 out physically adding more layers, making it efficient despite smaller model sizes. TroL
1865 uses CLIP-L and InternViT as vision encoders, containing 428M and 300M parameters,
1866 respectively, and supports 24 layers. The image resolution is adjusted using MLPs in the
1867 vision projector. For its foundational LLM, TroL utilizes Phi-3-mini with 3.8B parameters
1868 and InternLM2 with 1.8B and 7B parameters. The training dataset comprises 2.3M visual
1869 instruction tuning samples.
- 1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889