

BENCHMARKING VISUAL COGNITION OF MULTI-MODAL LLMs VIA MATRIX REASONING

Anonymous authors

Paper under double-blind review

ABSTRACT

Recently, Multimodal Large Language Models (MLLMs) and Vision Language Models (VLMs) have shown great promise in language-guided perceptual tasks such as recognition, segmentation, and object detection. However, their effectiveness in addressing visual cognition problems that require high-level multi-image reasoning and visual working memory is not well-established. One such challenge is matrix reasoning – the cognitive ability to discern relationships among patterns in a set of images and extrapolate to predict subsequent patterns. [This skill is crucial during the early neurodevelopmental stages of children.](#) Inspired by the matrix reasoning tasks in Raven’s Progressive Matrices (RPM) and Wechsler Intelligence Scale for Children (WISC), we propose a new dataset MaRs-VQA and a new benchmark VCog-Bench to evaluate the zero-shot visual cognition capability of MLLMs and compare their performance with existing [human visual cognition investigation](#). Our comparative experiments with different open-source and closed-source MLLMs on the VCog-Bench revealed a gap between MLLMs and human intelligence, highlighting the visual cognitive limitations of current MLLMs. We believe that the public release of VCog-Bench, consisting of MaRs-VQA, and the inference pipeline will drive progress toward the next generation of MLLMs with human-like visual cognition abilities.

1 INTRODUCTION

Matrix reasoning is a crucial ability in human perception and cognition, essential for nonverbal, culture-reduced intelligence measurements as it can minimize the influence of acquired knowledge and skills (Jensen, 1998; Jaeggi et al., 2010; Laurence & Macedo, 2023). Common matrix reasoning problems consist of images with simple shapes governed by underlying abstract rules (Małkiński & Mańdziuk, 2023) (see Figure 1). Participants have to identify and comprehend the rules based on a few provided patterns, and then reason about the next pattern following the same rules. Matrix reasoning is an important reflection of many fundamental capabilities of human intelligence, such as processing speed and working memory, that emerge in the early stage of children’s neurodevelopment (Gentner, 1977). To quantitatively measure human’s intelligence using matrix reasoning, many assessment methods have been proposed as a part of fluid intelligence tests. The two most famous assessments are Wechsler Intelligence Scale for Children (WISC) (Wechsler & Kodama, 1949) and Raven’s Progressive Matrices (RPM) (Raven, 2003).

In computer vision, matrix reasoning tasks have emerged as an ideal testbed for investigating whether deep learning models can match or even surpass human cognitive abilities, motivating the creation of diverse problem settings and datasets (Chollet, 2019; Małkiński & Mańdziuk, 2023; Barrett et al., 2018; Zhang et al., 2019; Webb et al., 2020). Previous research on matrix reasoning assessments applied typical machine learning settings – finetuning models on training sets and evaluating the performance on test sets (Hu et al., 2021; Małkiński & Mańdziuk, 2022; Zhao et al., 2024). [However, in human psychometrics, matrix reasoning are designed to assess visual reasoning abilities without prior specific training on similar tasks, which is similar to the zero-shot learning problem, but not training-testing paradigm in machine learning.](#) Children taking these tests typically do not receive any specialized training in matrix reasoning beforehand. Instead, they rely on their general cognitive skills developed through everyday experiences in natural scenes. Previous machine learning models ignore these prerequisites when modeling matrix reasoning problem. This could lead to an

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

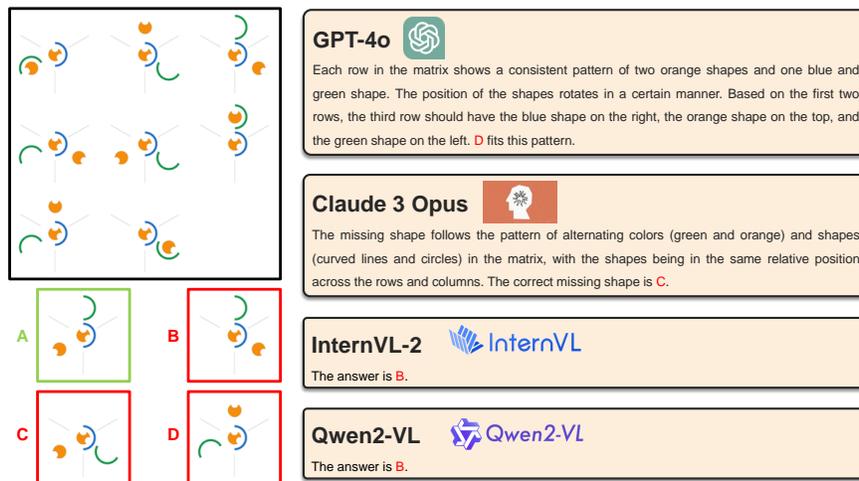


Figure 1: The example of the subpar performance of current state-of-the-art MLLMs (GPT-4o, Claude 3 Opus) and open-sourced VLMs (InternVL-2, Qwen2-VL) on a simple matrix reasoning task used in MaRs-VQA (similar to cases in RPM and WISC). Both models can recognize the basic shapes in the provided patterns but fail to reason the next pattern.

overestimation of the models’ reasoning abilities, as they might be overfitting with training data rather than demonstrating genuine generalization and reasoning skills from visual cognition.

Recently, Multimodal Large Language Models (MLLMs) have shown surprising understanding and reasoning capabilities, marking an important milestone towards Artificial General Intelligence (AGI) (Chollet, 2019; Ji et al., 2022; Peng et al., 2023). These models are learned from a large amount of data in the general domain and are proven can be generalize to unfamiliar tasks without prior exposure by in-context learning. However, current MLLMs remain inadequate in visual cognition problems that require higher-level inductive reasoning (Yang et al., 2023). An example is their poor performance on the RAVEN IQ-test (Huang et al., 2024; Fu et al., 2024), which heavily relies on abstract reasoning skills. The RAVEN IQ-test also has some limitations, including a small dataset of only 50 samples (Huang et al., 2024), which may introduce randomness and fail to comprehensively and robustly evaluate MLLMs. Besides, it doesn’t include a comparative study with human performance.

To address the matrix reasoning assessment and the deficiencies of existing cognitive testing benchmarks, we propose a new visual question answering (VQA) dataset – MaRs-VQA, which is the largest psychologist-verified VQA dataset for matrix reasoning assessment including 1,440 examples in total. The sample diversity of MaRs-VQA also surpasses other datasets before. It contains over 50 types of shape, 16 types of colour and over 500 graphic combinations. We also introduce VCog-Bench, the first zero-shot matrix reasoning benchmark to evaluate MLLMs’ visual cognition. In VCog-Bench, We conduct thorough evaluation and comparison across 16 existing MLLMs (including their variants) and human performance under a zero-shot inference setting (no prior knowledge) on MaRs-VQA and other abstract reasoning datasets containing human studies. In our experiments, we observe that MLLMs with more parameters generally perform better on our benchmark, adhering to established scaling laws in a limited scope. However, even the largest open-source MLLMs and GPT-4o fall short of surpassing human performance in matrix reasoning tasks. Furthermore, many MLLMs have a mismatch in performance between matrix reasoning tasks and other general VQA benchmarks, which provides some insights into the drawbacks of existing models. In conclusion, our contributions are summarized as follows:

- We introduce a new matrix reasoning VQA dataset – MaRs-VQA, containing 1,440 image instances designed by psychologists, which is the largest dataset for matrix reasoning zero-shot evaluation.
- We propose VCog-Bench, the most comprehensive visual cognition benchmark to date, which evaluates the matrix reasoning performance of 16 existing MLLMs and comparing them with human’s performance.

- Our thorough experiments qualitatively reveal the visual cognition gap between MLLMs and humans in matrix reasoning problems. We also show additional insights of deficiencies in MLLMs, which can inspire more future investigations in model design.

2 RELATED WORKS

Dataset	Source	Sample	Instance	RGB image	Human Study	Psychological Validity	Open-source	VQA Annotation
kosmos-iq50 (NeurIPS-23) (Huang et al., 2024)	RAVEN-IQ Test		50	✗	✗	✓	✗	✗
Visual Reasoning Benchmark (COLM-24) (Zhang et al., 2024c)	Mensa Test, RAVEN, IntelligenceTest		241	✗	✗	✗	✗	✗
MaRs-VQA (ours)	MaRs-IB		1,440	✓	✓	✓	✓	✓

Table 1: Comparison of recently released zero-shot matrix reasoning datasets to evaluate MLLMs.

Cognitive Test of Large Language Models (LLMs) The rise of LLMs has aroused interest in exploring human-like AI in psychology and cognition (Ullman, 2023). Recent works tested LLMs’ cognitive abilities in causal reasoning (Binz & Schulz, 2023), abstract reasoning (Xu et al., 2023b; Moskvichev et al., 2023; Jiang et al., 2024b; Ahrabian et al., 2024), analogical reasoning (Webb et al., 2023), systematic reasoning (Hagendorff et al., 2023), and theory of mind (Strachan et al., 2024). Their observation showed that LLMs like GPT-4 (Achiam et al., 2023) have been proven successful in most cognitive tests related to language-based reasoning. Despite this success, only limited research has been conducted on the areas of MLLMs and visual cognition. Visual cognition involves the process by which the human visual system interprets and makes inferences about a visual scene using partial information. *Buschoff et al.* observed that while LLMs demonstrate a basic understanding of physical laws and causal relationships, they lack deeper insights into intuitive human preferences and reasoning. Almost all existing visual cognition benchmarks focus on testing MLLMs’ cognitive abilities in simple tasks (Lerer et al., 2016; Zhou et al., 2023; Jassim et al., 2023), and ignore testing complex abstract reasoning and logical reasoning ability related to fluid intelligence. Therefore, new and challenging benchmarks based on the theory of visual cognition are needed to assess and improve AI systems’ capabilities for human-like visual understanding.

Matrix Reasoning Matrix reasoning is often used to determine human intelligence related to visual cognition and working memory (Salthouse, 1993; Jaeggi et al., 2010; Fleuret et al., 2011) that is widely used by RPM (Raven, 2003; Soulières et al., 2009), WISC (Wechsler & Kodama, 1949; Kaufman et al., 2015) to evaluate human’s ability to detect the underlying conceptual relationship among visual objects and use reasoning to find visual cues. Early research indicated that deep learning models can be trained with large-scale matrix reasoning datasets to solve simple matrix reasoning (Stabinger et al., 2021; Małkiński & Mańdziuk, 2022; 2023; Xu et al., 2023a; Małkiński & Mańdziuk, 2024) and compositional visual relation tasks (Fleuret et al., 2011; Zerroug et al., 2022; Ommer & Buhmann, 2007; Liu et al., 2021), achieving human-level accuracy. Several datasets and benchmarks are also proposed, such as PGM (Barrett et al., 2018), RAVEN (Zhang et al., 2019), RAVEN-I (Hu et al., 2021), RAVEN-FAIR (Benny et al., 2021), CVR (Zerroug et al., 2022). However, these works have a key limitation. They ignore that humans can solve these problems by zero-shot reasoning without explicitly learning from large-scale data. After the blooming of LLMs, researchers are keen on testing whether LLMs reached the same abstract reasoning capabilities as humans. *Webb et al.* (Webb et al., 2023) encode matrix reasoning into a symbolic problem based on human’s prior and validate LLM can understand this task. Recently, there are also some useful zero-shot visual reasoning inference datasets containing matrix reasoning samples have been proposed in the AI/ML community, such as RAVEN-IQ (Huang et al., 2024) containing 50 instances, Visual Reasoning Benchmark (Zhang et al., 2024c) containing 241 instances in total, but all of them are limited by

lacking rigorous human experiments as reference and conducting experiments on relatively small datasets without psychometrical validation.

Vision-Language Models Researchers have been actively investigating the utility of Vision-Language Models (VLMs) for addressing vision reasoning tasks (Zellers et al., 2019; Bordes et al., 2024). These latest VLMs are constructed using a combination of the CLIP vision encoder, pre-trained LLMs, and a connected adapter to align visual features with language space (Zhang et al., 2024b; Shao et al., 2024; Gupta & Kembhavi, 2023; Fu et al., 2024). Notably, methodologies such as MiniGPT-4 (Zhu et al., 2023), InstructBLIP (Dai et al., 2024), LLaVA (Liu et al., 2024b), CogVLM (Wang et al., 2023) underscore the significance of employing high-quality visual instruction tuning data. Additionally, tool learning methods have also explored the potential of integrating code generation pipelines with visual inference (Surís et al., 2023). Nevertheless, current VLMs encounter challenges in adapting to high-resolution and visually complex images. These problems stem from the absence of a robust visual search mechanism (Wu & Xie, 2023), few-shot reasoning (Guo et al., 2023), compositional understanding (Yuksekgonul et al., 2022) and the constrained visual grounding capabilities inherent in CLIP (Tong et al., 2024).

3 MARS-VQA DATASET

The MaRs-VQA dataset is designed to evaluate the zero-shot abstract reasoning capabilities of MLLMs through various matrix reasoning VQA tasks. The images in MaRs-VQA are sourced from the questionnaires in Matrix Reasoning Item Bank (MaRs-IB) (Chierchia et al., 2019), which is created by psychologists including 18 sets of abstract reasoning questionnaires (80 instances in each set) for non-verbal abstract reasoning assessment of adolescents and adults. Each item presents an incomplete 3×3 matrix of abstract shapes, requiring participants to identify relationships among the shapes. **Then, we create VQA annotations in the images from all questionnaires.**

To transform the matrix reasoning problem into a VQA task, we firstly define three different option sets – two image-based sets (A and B) and one language-based set (C). In Option Set A, we provide four candidates to the missing patch in the question. In Option Set B, the options are created by filling the four patches in Set A into the 3×3 question image. Note that Option Set B is used for visualization purposes only and is not included in our experiment. We further diversify the modalities of our dataset to support the evaluation of different kinds of models. Specifically, we use GPT-4o and human annotators to generate language-based descriptions for each option, forming Option Set C. In the data generation process, we first manually design 10 VQA examples, which serve as the initial human annotations in our data collection. These examples are then used as few-shot samples to query GPT-4o through in-context learning. The context generation system prompt guides GPT-4o to compare all four option images and generate distinct descriptions for each one. After generating all samples, human annotators in the author team review each option and revise the incorrect description. Examples are showed in Figure 6 in the Appendix. In Table 1, **Compared with other matrix reasoning datasets for MLLM’s visual cognition evaluation, MaRs-VQA is the largest one with unique features on psychological validity, human study reference, VQA annotations.**

4 VISUAL COGNITION BENCHMARK (VCOG-BENCH)

Different from the **training-testing paradigm** setting in other abstract visual reasoning datasets like RAVEN (Zhang et al., 2019), our goal of MLLM agent in VCog-Bench is to complete the 3×3 matrix by finding the missing cell from multiple options by **zero-shot learning under the same setting in human’s matrix reasoning test**. To this end, MLLM agents have to deduce relationships across the other cells of the matrix and infer the missing cell accordingly. **Based on the current progress of Multimodal LLMs, we propose two potential solutions as baselines for VCog-Bench.**

4.1 MULTI-IMAGE REASONING VIA CHAIN-OF-THOUGHT (COT)

Recent research in the NLP community has revealed the effectiveness of CoT in improving the reasoning capability of LLMs for complex problems (Wei et al., 2022; Kojima et al., 2022). In this paper, we propose the object-centric CoT prompting strategy, which combines the ideas of CoT (Zhang et al., 2023; Zhou et al., 2024; Zhang et al., 2024a), object-centric relational abstraction (Webb et al.,

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

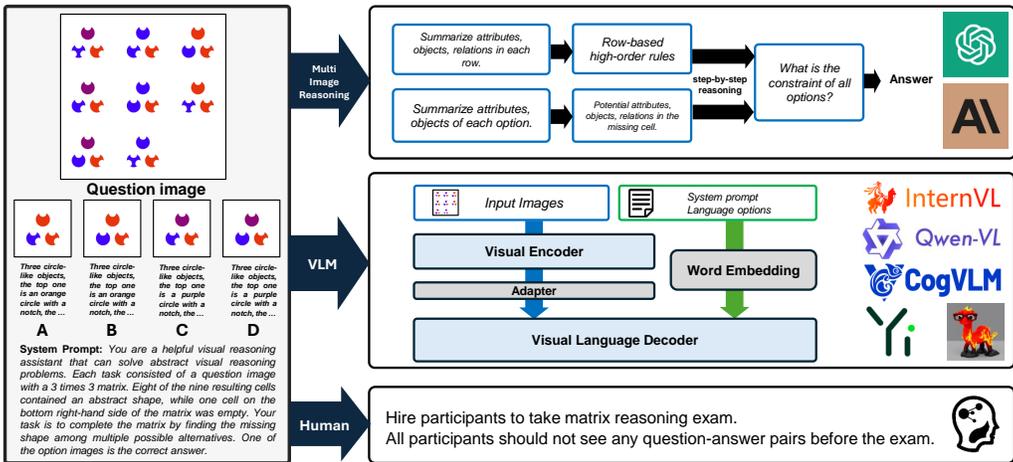


Figure 2: An overview of the VCog-Bench. The left part is the model input, including a question image, multiple option images and a system prompt describing the task. The right part shows the step-by-step CoT for multi-image reasoning and VLM solution for matrix reasoning problems.

2024a;b; Mondal et al., 2024; Xu et al., 2023b) and object-centric representation learning (Seitzer et al., 2022; Dittadi et al., 2022; Jiang et al., 2024a), to enhance the MLLM’s zero-shot learning performance in solving matrix reasoning problems.

Following previous works (Carpenter et al., 1990; Barrett et al., 2018; Chierchia et al., 2019; Zhang et al., 2019), we formulate the structure K of matrix reasoning as a combination of four components, $K = \{[r, a, o, s] | r \in \mathcal{R}, a \in \mathcal{A}, o \in \mathcal{O}, s \in \mathcal{S}\}$. \mathcal{R} is a set of rules of how the pattern changes along each row and column (e.g., rotating by a fixed angle and shifting by a fixed distance); \mathcal{A} is a set of attributes in each pattern (e.g., color, shape, and size); \mathcal{O} is how to integrate objects in each cell (e.g., spatial location and overlap); \mathcal{S} denotes a set of constraints for designing answer options (e.g., options should have minimum difference), which avoids that participants solving the matrix reasoning problems in unintended ways.

Based on structure K , we use three stages to guide MLLM to use human-level thought to understand matrix reasoning tasks. The first stage is to guide the Multimodal LLM to summarize the visual feature (e.g. shape) of each row in the 3×3 question image. Then, based on these row-based visual features, the model will then conclude the high-order rule/pattern \mathcal{R} . The second stage is to extract the basic attributes \mathcal{A} and inner relations \mathcal{O} to integrate objects in each option image. The third stage is to infer the answer based on exclusion with potential answer designed constraints \mathcal{S} . The system prompt of CoT will guide MLLM to step-by-step infer the sub-conclusion of each stage. And finally give the answer. The Multi-Image Reasoning section of Figure 2 shows a schematic depiction of how to leverage CoT in matrix reasoning tasks.

To further enhance CoT with diverse prompts, we introduce a multi-round architecture (Figure 3) inspired by the Monte Carlo Tree Search from Tree-of-Thought (ToT)(Yao et al., 2024). In the first reasoning round, the MLLM apply multi-image CoT solve the matrix reasoning problem. The selected image is then incorporated into the question image as a new input, which is fed back into the MLLM with a prompt directing it to evaluate the correctness of the complete 3×3 matrix, specifically focusing on the bottom-right corner. If the MLLM determines the result is correct, the final answer is output; otherwise, the incorrect option is excluded and CoT process is repeated.

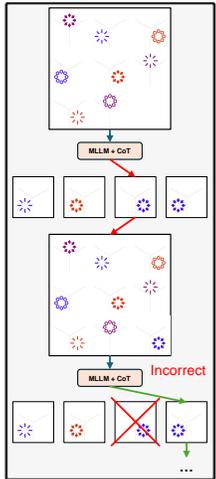


Figure 3: Multi-round CoT.

4.2 VISION-LANGUAGE MODELS (VLMs)

In addition to MLLMs, we also evaluate the performance of VLMs for a thorough comparison. In VLMs, we only use question image as visual input and transform all option images into language

descriptions (*i.e.*, Option Set C), which matches the input representations required by VLMs (Xu et al., 2023b; Camposampiero et al., 2023). The VLM section in Figure 2 illustrates this pipeline.

The test set contains n VQA samples, denoted as $\{(\mathbf{q}_i, \mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$. \mathbf{q}_i represents the question image showing the 3×3 matrix reasoning task (MaRs-VQA). $\mathbf{x}_i = [x_i^1, \dots, x_i^k]$ represents the context description in the option set, where k is the number of options. \mathbf{y}_i is the answer of the matrix reasoning question. The zero-shot inference pipeline of VLM can be formulated as:

$$\hat{\mathbf{y}}_i = F_\theta(\mathbf{q}_i, \mathbf{x}_i, \mathbf{x}_{sys}). \quad (1)$$

\mathbf{x}_{sys} is the system prompt, including independent information about the matrix reasoning problem setting, structure K for each dataset and requirements for the output format. $\hat{\mathbf{y}}_i$ is the prediction result. F_θ is an autoregressive decoder in the LLM for answer generation. It is defined as:

$$P(\hat{\mathbf{y}}_i | \mathbf{q}_i, \mathbf{x}_i, \mathbf{x}_{sys}) = \prod_{j=1}^L P(\hat{\mathbf{y}}_{i,j} | f(\mathbf{q}_i), \mathbf{x}_i, \mathbf{x}_{sys}, \hat{\mathbf{y}}_{i,<j}; \theta), \quad (2)$$

where f is the visual encoder and adapter layer, L is the sequence length of answers and $\hat{\mathbf{y}}_{i,<j}$ is all answer tokens before $\hat{\mathbf{y}}_{i,j}$.

In VLMs, the input question image is first processed by the visual encoder such as CLIP (Radford et al., 2021). Then, additional adapter layers are used to map visual features into language feature space. These features, along with the context-based option descriptions, are sent to the LLM decoder. The LLM decoder then integrates the information from both the input question image and the option descriptions to address the VQA task. VLMs leverage the strengths of both visual encoders and language models, allowing for a more comprehensive analysis of the matrix reasoning problems. It provides a structured way to break down the problem, potentially improving interpretability compared to end-to-end close-source models.

5 EXPERIMENTS

5.1 EXPERIMENTAL SETTINGS

Datasets In addition to MaRs-VQA, we selected two well-known open-source datasets for matrix reasoning and abstract visual reasoning to conduct experiments in VCog-Bench. The first dataset is RAVEN (Zhang et al., 2019), designed to probe abstract reasoning in a format similar to the Raven’s Progressive Matrices IQ test, with each question providing eight options. The second dataset is Compositional Visual Reasoning (CVR) (Zerroug et al., 2022), which evaluates deep learning models using 103 unique configurations generated by predefined rules. Each sample in CVR is an outlier detection problem, with four options provided per question. However, both RAVEN and CVR share a significant limitation: all samples are algorithmically generated using fixed rules, which limits their diversity and lacks psychological validity.

Baselines for Multi-image Reasoning We selected the Claude 3 family (Haiku, Sonnet, Opus) (Anthropic, 2024), GPT-4V (OpenAI, 2023), GPT-4o (OpenAI, 2024) as the primary multi-image CoT baselines as they support multiple images input and can generate reasoning process. [The inputs for this task are all images, a question and multiple option images in Option Set A of Figure 6. Other open-sourced models are not included because they perform much worse than Claude and GPT and can not generate reasoning steps for matrix reasoning tasks.](#)

Baselines for VLMs For the VLMs, we select state-of-the-arts open-source and closed-source models such as InstructBLIP (Dai et al., 2024), MiniGPT-v2 (Zhu et al., 2023), LLaVA-v1.6 (LLaVA-NeXT) (Liu et al., 2024a), CogVLMv2 (Wang et al., 2023), Yi-VL (Young et al., 2024), Qwen-VL (Bai et al., 2023), InternVL (Chen et al., 2024), Gemini Pro 1.5 (Reid et al., 2024), Claude 3 family (Haiku, Sonnet, Opus) (Anthropic, 2024), GPT-4V (OpenAI, 2023), GPT-4o (OpenAI, 2024) as the primary VLM baselines. The input is a question image and language-based options.

Human Baseline The human study results in Table 2 and 3 are reported from previous experiment results. The human subjects of RAVEN (Zhang et al., 2019) consists of college students from a

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

Method	Learning	Accuracy (%) \uparrow		
		MaRs-VQA (4-options)	RAVEN (8-options)	CVR (4-options)
Claude 3 Sonnet (Anthropic, 2024)	zero-shot	22.92	10.71	27.83
	CoT	23.22	13.39	28.48
Claude 3 Opus (Anthropic, 2024)	zero-shot	20.85	11.61	26.86
	CoT	24.13	11.95	27.18
Claude 3.5 Sonnet (Anthropic, 2024)	zero-shot	23.18	14.08	25.97
	CoT	24.28	15.36	27.88
GPT-4V (OpenAI, 2023)	zero-shot	27.71	13.84	36.25
	CoT	33.13	15.63	40.62
GPT-4o (OpenAI, 2024)	zero-shot	30.21	19.20	42.50
	CoT	33.96	25.89	44.01
Human	-	69.15	84.41	78.70

Table 2: Experiments on multi-image reasoning. zero-shot means only provide the model system prompt about the matrix reasoning task definition. Chain-of-thought denotes the implementation in section 4.1. The results are averaged over three runs with three different random seeds.

Method	Training Data	Model Scale	LLM Backbone	Accuracy (%) \uparrow	
				MaRs-VQA (4 Options)	RAVEN (8 Options)
InstructBLIP (Dai et al., 2024)	129M	7B	Vicuna-7B (Chiang et al., 2023)	10.63	12.05
LLaVA-v1.6 (Liu et al., 2024b)	1.3M	7B	Mistral-7B (Jiang et al., 2023)	16.88	14.29
MiniGPT-v2 (Zhu et al., 2023)	-	8B	Llama-2-7B (Touvron et al., 2023)	26.45	13.39
Qwen-VL (Bai et al., 2023)	1.4B	10B	Qwen-7B (Bai et al., 2023)	29.58	16.07
InstructBLIP (Dai et al., 2024)	129M	13B	Vicuna-13B (Chiang et al., 2023)	10.42	14.46
CogVLMv2 (Wang et al., 2023)	1.5B	19B	Llama-3-8B (Meta, 2024a)	26.46	12.05
InternVL 1.5 (Chen et al., 2024)	6.0B	26B	InternLM2-Chat-20B (Cai et al., 2024)	22.09	14.73
Yi-VL (Young et al., 2024)	100M	34B	Yi-34B-Chat (Young et al., 2024)	25.21	19.64
LLaVA-v1.6 (Liu et al., 2024b)	1.3M	35B	Hermes-Yi-34B (Young et al., 2024)	34.38	33.93
InternVL 1.2+ (Chen et al., 2024)	6.0B	40B	Hermes-Yi-34B (Young et al., 2024)	32.71	33.04
Qwen2-VL (Wang et al., 2024)	-	72B	Qwen2-72B (Yang et al., 2024)	34.22	36.15
InternVL 2 (Chen et al., 2024)	-	76B	Hermes-2-Theta-Llama-3-70B (Teknium et al.)	34.63	38.01
Llama 3.2 (Meta, 2024b)	6.0B	90B	-	34.81	35.26
Claude 3.5 Sonnet (Anthropic, 2024)	<i>unknown</i>	<i>unknown</i>	<i>unknown</i>	34.82	35.36
GPT-4o (OpenAI, 2024)	<i>unknown</i>	<i>unknown</i>	<i>unknown</i>	37.38	38.84
Gemini Pro 1.5 (Reid et al., 2024)	<i>unknown</i>	<i>unknown</i>	<i>unknown</i>	34.79	42.86
Human	-	-	-	69.15	84.41

Table 3: Experiments on using a question image and language descriptions for options as inputs to compare different VLMs. The results are averaged over three random seeds.

subject pool maintained by the Department of Psychology. Only “easily perceptible” examples were used in the investigation. CVR (Zerroug et al., 2022) hired 21 participants and each participant completed 6 different tasks with 20 problem samples for each task. The human study results of MaRs-IB (Chierchia et al., 2019) (data source of MaRs-VQA) are more rigorous. They are from 4 age groups ($N = 659$, aged 11–33 years). The accuracy for younger adolescents, mid-adolescents, older adolescents, and adults solving matrix reasoning in MaRs-IB are 61%, 68%, 73%, 81%. We use the average result of all groups in Table 2 and 3.

Implementation For closed-source baseline models, we establish basic prompts to introduce the matrix reasoning problem setting, which serve as the system prompt for zero-shot inference. For object-centric CoT reasoning, we create specific prompts to guide the model’s thought process through multiple stages, enabling step-by-step reasoning. For open-source baseline models, we use the same system prompt settings across all models. Testing is conducted using two NVIDIA RTX 4090 GPUs for 7B-sized VLMs and eight NVIDIA A100 80GB GPUs for VLMs larger than 7B. All experiments are run with three different random seeds, and the results are averaged. We evaluate the results based on the accuracy of single-option matrix reasoning problems ($\text{Acc} = \text{Correct}/\text{Total}$), consistent with other VQA benchmarks (Lu et al., 2022; Liu et al., 2023).

5.2 EXPERIMENTAL RESULTS

In this subsection, we present the experimental results of the baselines in the VCog-Bench. The results demonstrate that while parts of baseline models can understand some basic forms of the matrix reasoning task, they struggle with complex tasks requiring both visual working memory and multi-image reasoning capability.

We divided our experiments into two parts. The first part involves end-to-end multi-image reasoning. For this experiment, we used multiple images as the input, including a question image and several option images (refer to Option Set A in Figure 6), and guided the MLLMs to decompose the problem into predefined structures before generating answers based on all available information. We tested

Method	Multi-Image	Accuracy (%) \uparrow				
		Level 1 (90)	Level 2 (96)	Level 3 (84)	Level 4 (72)	Level >4 (138)
Claude 3 Opus (Anthropic, 2024)	✓	19.15	28.57	13.34	13.16	24.66
GPT-4o (OpenAI, 2024)	✓	57.78	27.08	27.38	19.43	21.74
Claude 3 Opus (Anthropic, 2024)	✗	24.44	25.00	40.48	38.89	39.13
Gemini Pro 1.5 (Reid et al., 2024)	✗	51.10	30.21	26.19	29.17	35.51
GPT-4o (OpenAI, 2024)	✗	58.89	45.83	32.14	26.39	26.09

Table 4: Compare closed-source MLLMs with different difficulty levels in MaRs-VQA. The number in the “()” is the number of case sample of selected level. The difficulty level is based on the complexity of color, size, geometry, positional relationships, and object counting.

the Claude 3 family, GPT-4V, and GPT-4o for this task, as these models can generate step-by-step multi-image reasoning. Table 2 shows that even the state-of-the-art closed-source MLLMs perform worse than humans in all matrix reasoning tasks. While object-centric CoT can help larger models achieve better performance, it does not benefit smaller models such as Claude 3 Sonnet. Compared to the results in MaRs-VQA and RAVEN, GPT-4o achieves much better zero-shot and object-centric CoT inference results in the CVR dataset, almost matching the performance (ResNet-50: 57.9%, ViT-small: 32.7%, WReN: 42.4%) of fine-tuned models with 1,000 training samples in CVR’s paper (Zerroug et al., 2022).

In the second part of our experiment, we investigated the use of VLMs (question image + language options) to solve matrix reasoning problems in MaRs-VQA and RAVEN. The CVR dataset was excluded because the shapes it contains are too complex to describe accurately. As shown in Table 3, large-scale VLMs, such as Qwen2-72B and InternVL-2-76B, achieved comparable results to GPT-4o in MaRs-VQA and RAVEN. Notably, Gemini Pro 1.5 outperformed GPT-4o on the RAVEN dataset.

We identified three major issues after reviewing the reasoning outputs of current MLLMs in Table 2 and 3: (1) Limited Use of Visual Information: MLLMs cannot directly use visual features for reasoning, making them insensitive to non-verbal spatial features during CoT reasoning. This limitation is particularly evident when handling images that require describing the positional relations of objects. For example, it is difficult for MLLMs to distinguish each option in Figure 1 using language alone. (2) Restricted Visual Working Memory: The visual working memory of MLLMs is limited, causing visual feature information to be easily lost during the text generation reasoning process. (3) Integration Challenges: Even if MLLMs possess strong task-specific skills like recognition, segmentation, and object detection, they struggle to integrate these skills into high-level visual reasoning tasks. We will further analyse them in the ablation study.

5.3 ABLATION STUDY

In this subsection, we conduct ablation experiments to analyze how to improve the performance of MLLMs on the matrix reasoning problem. Table 5 compares the Chain-of-Thought (CoT) baseline with two approaches: few-shot reasoning and multi-round reasoning. Few-shot reasoning involves providing a small number of question-answer examples alongside the CoT system prompt. Multi-round reasoning employs the advanced CoT strategy illustrated in Figure 3. The results show that incorporating 1-shot and 3-shot question-option-answer pairs gradually increases the accuracy on MaRs-VQA from 34% to 36%. However, extending the number of examples to 5 does not yield further improvement. These findings suggest that while few-shot in-context learning helps the model better understand the matrix reasoning problem, it does not significantly enhance the MLLM’s visual reasoning capabilities for these tasks. Additionally, using a multi-round tree search improves accuracy from approximately 34% to 42%, but it is considerably slower than single-round CoT, with each inference taking over 30 seconds in multi-round mode. We also compare different MLLMs across difficulty levels and different visual complexities in the MaRs-VQA dataset (see Table 4 and Table 6). The difficulty level in our tasks is defined by the number of sub-tasks (visual complexities) involved—specifically, variations in color, size, geometry (shape), positional relationships, and the presence of multiple objects. The results indicate that GPT-4o exhibits difficulty sensitivity similar to that of humans, whereas Claude 3 Opus does not demonstrate this ability. This is because GPT-4o can solve object size sub-tasks well in the MaRs-VQA, but is still struggling with other sub-tasks, especially the multi-object sub-task. More details are presented in the Appendix.

Strategy	Accuracy (%) \uparrow
CoT	33.96
CoT + 1-shot	35.22
CoT + 3-shot	36.10
CoT + 5-shot	36.03
multi-round tree search CoT	41.96
multi-round tree search CoT + 1-shot	42.08

Table 5: Ablation on prompt selection.

Visual complexity	Proportion (%)	Accuracy (%) \uparrow
Shape	68	33.96
Color	73	35.72
Size	16	63.26
Position	41	31.70
Multi-Object	71	31.48
All	-	33.96

Table 6: Ablation on visual complexity.

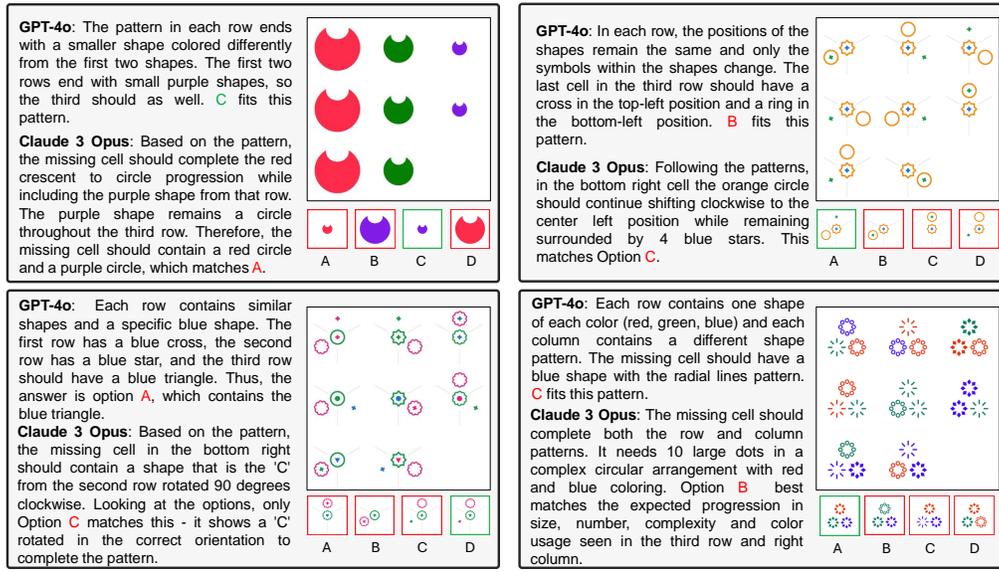


Figure 4: Different matrix reasoning problem (difficulty levels) from MaRs-VQA and MLLM’s reply. We use green to represent correct answer and red to represent wrong answer of each question. The top left is a sample with difficulty level 1. The others are samples with difficulty level ≥ 4 . The reasoning is a short summary of the CoT output, not the full version

5.4 QUALITATIVE ANALYSIS

In this subsection, we use case studies from the MaRs-VQA dataset to illustrate how MLLMs fail in some tasks and provide insights on how to improve MLLMs and VLMs for this task.

First, we present an example to explain why the Claude 3 family performs worse than GPT-4o and even worse than random guessing in most of experiments. Figure 4 top left is one of the most simple cases in MaRs-VQA’s level 1 difficulty, Claude 3 Opus incorrectly identifies the shape as the main target of this matrix, while the actual target is the size. In contrast, GPT-4o correctly discerns the relationship between rows, noting: “The pattern in each row ends with a smaller shape colored differently from the first two shapes.” This example highlights a critical shortcoming in Claude 3 Opus’s reasoning ability: limited use of Visual information, demonstrating its struggle to accurately interpret the key attributes in matrix reasoning tasks. GPT-4o, on the other hand, showcases a superior understanding of the relationships within simple data, leading to more accurate responses.

However, the difficulty of the tasks increases, the performance of MLLMs deteriorates in multi-image reasoning. Figure 4 bottom left and shows an example, it is the level 6 difficulty containing shape, positional relation, shape with different objects. For these questions containing complex visual features, MLLMs tend to extract only a small portion of the key information from the question image. This limited extraction means that critical features are either overlooked or not effectively utilized in selecting the correct option. Consequently, the final answers are often incorrect or only partially related to the relevant attributes. It suggests that MLLMs are affected by the cognitive load associated with processing multiple sub-tasks simultaneously, which is closely related to the concept of visual working memory. The right two examples of Figure 4 also present the same observation. Additionally, we observed that GPT-4o is not sensitive to the positional relationships for multi-objects in the question images.

These failures highlight significant limitations in MLLM’s visual processing capabilities. The model’s inability to effectively leverage visual features and its lack of visual working memory result in incorrect interpretations. Furthermore, its insensitivity to positional relationships among multi-objects underscores a critical area for improvement in understanding and analyzing spatial information in visual reasoning.

5.5 VISUALIZATION

We also analyze the relationship between matrix reasoning accuracy and model scale in Figure 5. The figure illustrates the significant gap between MLLM’s matrix reasoning performance and that of humans. This gap is substantial and suggests that simply increasing model size according to scaling laws will not be sufficient to bridge it.

6 DISCUSSION

Social Impacts In the present work, we emphasize that zero-shot matrix reasoning is a key item to validate human-level intelligence, though it is still unclear how matrix reasoning ability is acquired early in human neurodevelopment. Children’s visual reasoners (without any additional training) can provide sensible answers to matrix reasoning questions as early as age four. The long-term goal of our work is twofold. The first one is to explore the problem of how close AIs or MLLMs are to human-like cognitive abilities, which is raised by *François Chollet* in 2019 (Chollet, 2019). The second one is to develop an MLLM-powered AI agent that can simulate human-level zero-shot matrix reasoning capability. The agent will eventually guide vision generation models to generate new matrix reasoning samples and tasks and design new neurodevelopmental assessment tools. This will help psychologists and pediatricians explore and deconstruct how children activate such abilities in the early stage of neurodevelopment.

Limitations An open-ended question is whether MLLMs need to achieve or surpass human-level zero-shot inference capability in matrix reasoning tasks. Addressing this issue requires new theories from cognitive science and psychology to accurately evaluate and compare human and MLLM intelligence. Unlike MLLMs, which rely on training data and domain-specific skills, human cognition develops gradually and evolves with age. Humans can also learn how to solve the problem progressively from previously seen matrix reasoning tasks while they are taking the test, but MLLM can not learn from it via in-context learning due to the maximum tokens length. Therefore, AI researchers, psychologists, and cognitive scientists must collaborate to rethink how to benchmark MLLM intelligence with human intelligence.

7 CONCLUSION

We introduce VCog-Bench, a publicly available zero-shot matrix reasoning benchmark designed to evaluate the visual cognition capability and intelligence of Multimodal Large Language Models (MLLMs). This benchmark integrates two well-known datasets RAVEN and CVR from the AI community and includes our newly proposed MaRs-VQA dataset. We also introduce several important concepts to redefine zero-shot matrix reasoning task evaluation, focusing on multi-image reasoning with object-centric Chain-of-Thought (CoT) system prompts. Our findings show that current state-of-the-art MLLMs and Vision-Language Models (VLMs), such as GPT-4o and LLaVA-1.6, InternVL demonstrate some basic understanding of matrix reasoning tasks. However, these models still face big challenges with complex situations and perform much worse than human. This highlights the need for further exploration and development in this area. By providing a robust benchmark, we aim to encourage further innovation and progress in the field of improving the visual cognition of MLLMs.

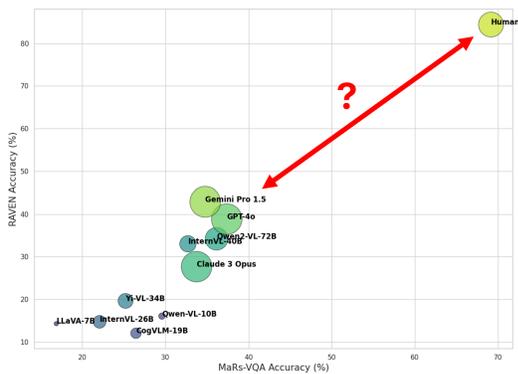


Figure 5: There is still a big gap between human’s matrix reasoning capability and MLLM’s. Bubble size corresponds to the model size. As we don’t know the exact size of closed-source MLLMs, we set all of them to the largest value by default. The model size of human refers to the number of neurons (86B) in human’s brain (Voytek, 2013).

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
543 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
544 *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Kian Ahrabian, Zhivar Sourati, Kexuan Sun, Jiarui Zhang, Yifan Jiang, Fred Morstatter, and Jay
546 Pujara. The curious case of nonverbal abstract reasoning with multi-modal large language models.
547 *arXiv preprint arXiv:2401.12117*, 2024.
- 548 Anthropic. Introducing the next generation of claude. [https://www.anthropic.com/news/
549 claude-3-family](https://www.anthropic.com/news/claude-3-family), 2024.
550
- 551 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang
552 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities.
553 *arXiv preprint arXiv:2308.12966*, 2023.
- 554 David Barrett, Felix Hill, Adam Santoro, Ari Morcos, and Timothy Lillicrap. Measuring abstract
555 reasoning in neural networks. In *International conference on machine learning*, pp. 511–520.
556 PMLR, 2018.
- 557 Yaniv Benny, Niv Pekar, and Lior Wolf. Scale-localized abstract reasoning. In *Proceedings of the
558 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12557–12565, 2021.
559
- 560 Marcel Binz and Eric Schulz. Using cognitive psychology to understand gpt-3. *Proceedings of the
561 National Academy of Sciences*, 120(6):e2218523120, 2023.
- 562 Florian Bordes, Richard Yuanzhe Pang, Anurag Ajay, Alexander C. Li, Adrien Bardes, Suzanne
563 Petryk, Oscar Mañas, Zhiqiu Lin, Anas Mahmoud, Bargav Jayaraman, et al. An introduction to
564 vision-language modeling. *arXiv preprint arXiv:2405.17247*, 2024.
565
- 566 Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen,
567 Zhi Chen, Pei Chu, et al. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*, 2024.
- 568 Giacomo Camposampiero, Loïc Houmar, Benjamin Estermann, Joël Mathys, and Roger Wattenhofer.
569 Abstract visual reasoning enabled by language. In *Proceedings of the IEEE/CVF Conference on
570 Computer Vision and Pattern Recognition*, pp. 2642–2646, 2023.
- 571 Patricia A Carpenter, Marcel A Just, and Peter Shell. What one intelligence test measures: a
572 theoretical account of the processing in the raven progressive matrices test. *Psychological review*,
573 97(3):404, 1990.
- 574 Raymond Bernard Cattell and Alberta KS Cattell. *Measuring intelligence with the culture fair tests*.
575 Institute for Personality and Ability Testing, 1960.
- 576 Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi
577 Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial
578 multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024.
- 579 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,
580 Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot
581 impressing gpt-4 with 90%* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April
582 2023), 2(3):6, 2023.
- 583 Gabriele Chierchia, Delia Fuhrmann, Lisa J Knoll, Blanca Piera Pi-Sunyer, Ashok L Sakhardande,
584 and Sarah-Jayne Blakemore. The matrix reasoning item bank (mars-ib): novel, open-access
585 abstract reasoning items for adolescents and adults. *Royal Society open science*, 6(10):190232,
586 2019.
- 587 François Chollet. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.
588
- 589 Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang,
590 Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-
591 language models with instruction tuning. *Advances in Neural Information Processing Systems*, 36,
592 2024.
593

- 594 Andrea Dittadi, Samuele S Papa, Michele De Vita, Bernhard Schölkopf, Ole Winther, and Francesco
595 Locatello. Generalization and robustness implications in object-centric learning. In *International*
596 *Conference on Machine Learning*, pp. 5221–5285. PMLR, 2022.
- 597
- 598 François Fleuret, Ting Li, Charles Dubout, Emma K Wampler, Steven Yantis, and Donald Geman.
599 Comparing machines and humans on a visual categorization test. *Proceedings of the National*
600 *Academy of Sciences*, 108(43):17621–17625, 2011.
- 601 Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith,
602 Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not
603 perceive. *arXiv preprint arXiv:2404.12390*, 2024.
- 604
- 605 Dedre Gentner. Children’s performance on a spatial analogies task. *Child development*, pp. 1034–
606 1039, 1977.
- 607
- 608 Qing Guo, Prashan Wanigasekara, Jian Zheng, Jacob Zhiyuan Fang, Xinwei Deng, and Chenyang
609 Tao. How do large multimodal models really fare in classical vision few-shot challenges? a deep
610 dive. In *R0-FoMo: Robustness of Few-shot and Zero-shot Learning in Large Foundation Models*,
611 2023.
- 612 Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning
613 without training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
614 *Recognition*, pp. 14953–14962, 2023.
- 615 Thilo Hagendorff, Sarah Fabi, and Michal Kosinski. Human-like intuitive behavior and reasoning
616 biases emerged in large language models but disappeared in chatgpt. *Nature Computational*
617 *Science*, 3(10):833–838, 2023.
- 618
- 619 Sheng Hu, Yuqing Ma, Xianglong Liu, Yanlu Wei, and Shihao Bai. Stratified rule-aware network
620 for abstract visual reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
621 volume 35, pp. 1567–1574, 2021.
- 622 Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv,
623 Lei Cui, Owais Khan Mohammed, Barun Patra, et al. Language is not all you need: Aligning
624 perception with language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- 625
- 626 Susanne M Jaeggi, Barbara Studer-Luethi, Martin Buschkuehl, Yi-Fen Su, John Jonides, and Walter J
627 Perrig. The relationship between n-back performance and matrix reasoning—implications for
628 training and transfer. *Intelligence*, 38(6):625–635, 2010.
- 629 Serwan Jassim, Mario Holubar, Annika Richter, Cornelius Wolff, Xenia Ohmer, and Elia Bruni.
630 Grasp: A novel benchmark for evaluating language grounding and situated physics understanding
631 in multimodal language models. *arXiv preprint arXiv:2311.09048*, 2023.
- 632
- 633 Arthur R Jensen. *The factor*. Westport, CT: Prager, 1998.
- 634
- 635 Anya Ji, Noriyuki Kojima, Noah Rush, Alane Suhr, Wai Keen Vong, Robert D Hawkins, and Yoav
636 Artzi. Abstract visual reasoning with tangram shapes. *arXiv preprint arXiv:2211.16492*, 2022.
- 637
- 638 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
639 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
640 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 641 Jindong Jiang, Fei Deng, Gautam Singh, and Sungjin Ahn. Object-centric slot diffusion. *Advances in*
642 *Neural Information Processing Systems*, 36, 2024a.
- 643 Yifan Jiang, Jiarui Zhang, Kexuan Sun, Zhivar Sourati, Kian Ahrabian, Kaixin Ma, Filip Ilievski, and
644 Jay Pujara. Marvel: Multidimensional abstraction and reasoning through visual evaluation and
645 learning. *arXiv preprint arXiv:2404.13591*, 2024b.
- 646
- 647 Alan S Kaufman, Susan Engi Raiford, and Diane L Coalson. *Intelligent testing with the WISC-V*.
John Wiley & Sons, 2015.

- 648 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
649 language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:
650 22199–22213, 2022.
- 651 Paulo Guirro Laurence and Elizeu Coutinho Macedo. Cognitive strategies in matrix-reasoning tasks:
652 State of the art. *Psychonomic Bulletin & Review*, 30(1):147–159, 2023.
- 653 Adam Lerer, Sam Gross, and Rob Fergus. Learning physical intuition of block towers by example.
654 In *International conference on machine learning*, pp. 430–438. PMLR, 2016.
- 655 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.
656 Llava-next: Improved reasoning, ocr, and world knowledge. <https://llava-vl.github.io/blog/2024-01-30-llava-next>, 2024a.
- 657 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in
658 neural information processing systems*, 36, 2024b.
- 659 Nan Liu, Shuang Li, Yilun Du, Josh Tenenbaum, and Antonio Torralba. Learning to compose visual
660 relations. *Advances in Neural Information Processing Systems*, 34:23166–23178, 2021.
- 661 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi
662 Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player?
663 *arXiv preprint arXiv:2307.06281*, 2023.
- 664 Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,
665 Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for
666 science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521,
667 2022.
- 668 Mikolaj Małkiński and Jacek Mańdziuk. Deep learning methods for abstract visual reasoning: A
669 survey on raven’s progressive matrices. *arXiv preprint arXiv:2201.12382*, 2022.
- 670 Mikolaj Małkiński and Jacek Mańdziuk. A review of emerging research directions in abstract visual
671 reasoning. *Information Fusion*, 91:713–736, 2023.
- 672 Mikolaj Małkiński and Jacek Mańdziuk. One self-configurable model to solve many abstract visual
673 reasoning problems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38,
674 pp. 14297–14305, 2024.
- 675 AI Meta. Introducing meta llama 3: The most capable openly available llm to date, 2024. URL
676 <https://ai.meta.com/blog/meta-llama-3/>. Accessed on April, 26, 2024a.
- 677 AI Meta. Introducing llama 3.2. URL https://github.com/meta-llama/llama-models/tree/main/models/llama3_2 Accessed on Sep, 2024b.
- 678 Shanka Subhra Mondal, Jonathan D Cohen, and Taylor W Webb. Slot abstractors: Toward scalable
679 abstract visual reasoning. *arXiv preprint arXiv:2403.03458*, 2024.
- 680 Arsenii Kirillovich Moskvichev, Victor Vikram Odouard, and Melanie Mitchell. The conceptarc
681 benchmark: Evaluating understanding and generalization in the arc domain. *Transactions on
682 Machine Learning Research*, 2023.
- 683 Bjorn Ommer and Joachim M Buhmann. Learning the compositional nature of visual objects. In
684 *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–8. IEEE, 2007.
- 685 OpenAI. Gpt-4v(ision) system card. [https://openai.com/research/
686 gpt-4v-system-card](https://openai.com/research/gpt-4v-system-card), 2023.
- 687 OpenAI. Hello gpt-4o. <https://openai.com/index/hello-gpt-4o>, 2024.
- 688 Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu
689 Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint
690 arXiv:2306.14824*, 2023.
- 691

- 702 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
703 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
704 models from natural language supervision. In *International conference on machine learning*, pp.
705 8748–8763. PMLR, 2021.
- 706
707 Jean Raven. Raven progressive matrices. In *Handbook of nonverbal assessment*, pp. 223–237.
708 Springer, 2003.
- 709
710 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste
711 Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini
712 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint*
713 *arXiv:2403.05530*, 2024.
- 714
715 Timothy A Salthouse. Influence of working memory on adult age differences in matrix reasoning.
716 *British Journal of Psychology*, 84(2):171–199, 1993.
- 717
718 Maximilian Seitzer, Max Horn, Andrii Zadaianchuk, Dominik Zietlow, Tianjun Xiao, Carl-Johann
719 Simon-Gabriel, Tong He, Zheng Zhang, Bernhard Schölkopf, Thomas Brox, et al. Bridging the
720 gap to real-world object-centric learning. In *The Eleventh International Conference on Learning*
Representations, 2022.
- 721
722 Hao Shao, Shengju Qian, Han Xiao, Guanglu Song, Zhuofan Zong, Letian Wang, Yu Liu, and
723 Hongsheng Li. Visual cot: Unleashing chain-of-thought reasoning in multi-modal language
724 models. *arXiv preprint arXiv:2403.16999*, 2024.
- 725
726 Isabelle Soulières, Michelle Dawson, Fabienne Samson, Elise B Barbeau, Cherif P Sahyoun, Gary E
727 Strangman, Thomas A Zeffiro, and Laurent Mottron. Enhanced visual processing contributes to
matrix reasoning in autism. *Human brain mapping*, 30(12):4082–4107, 2009.
- 728
729 Sebastian Stabinger, David Peer, Justus Piater, and Antonio Rodríguez-Sánchez. Evaluating the
730 progress of deep learning for visual relational concepts. *Journal of Vision*, 21(11):8–8, 2021.
- 731
732 James WA Strachan, Dalila Albergo, Giulia Borghini, Oriana Pansardi, Eugenio Scaliti, Saurabh
733 Gupta, Krati Saxena, Alessandro Rufo, Stefano Panzeri, Guido Manzi, et al. Testing theory of
mind in large language models and humans. *Nature Human Behaviour*, pp. 1–11, 2024.
- 734
735 Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for
736 reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.
737 11888–11898, 2023.
- 738
739 Teknium, Charles Goddard, interstellarninja, theemozilla, karan4d, and huemin_art.
740 Hermes-2-theta-llama-3-70b. [https://huggingface.co/NousResearch/
Hermes-2-Theta-Llama-3-70B](https://huggingface.co/NousResearch/Hermes-2-Theta-Llama-3-70B).
- 741
742 Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide
743 shut? exploring the visual shortcomings of multimodal llms. *IEEE/CVF Conference on Computer*
744 *Vision and Pattern Recognition (CVPR)*, 2024.
- 745
746 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
747 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
748 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 749
750 Tomer Ullman. Large language models fail on trivial alterations to theory-of-mind tasks. *arXiv*
preprint arXiv:2302.08399, 2023.
- 751
752 Bradley Voytek. Are there really as many neurons in the human brain as stars in the milky way.
753 *Scitable, Nature Education*, 2013.
- 754
755 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,
Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the
world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.

- 756 Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang,
757 Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv*
758 *preprint arXiv:2311.03079*, 2023.
- 759 Taylor Webb, Keith J Holyoak, and Hongjing Lu. Emergent analogical reasoning in large language
760 models. *Nature Human Behaviour*, 7(9):1526–1541, 2023.
- 761 Taylor Webb, Shanka Subhra Mondal, and Jonathan D Cohen. Systematic visual reasoning through
762 object-centric relational abstraction. *Advances in Neural Information Processing Systems*, 36,
763 2024a.
- 764 Taylor W Webb, Steven M Frankland, Awni Altabaa, Simon Segert, Kamesh Krishnamurthy, Declan
765 Campbell, Jacob Russin, Tyler Giallanza, Randall O’Reilly, John Lafferty, et al. The relational
766 bottleneck as an inductive bias for efficient abstraction. *Trends in Cognitive Sciences*, 2024b.
- 767 Taylor Whittington Webb, Ishan Sinha, and Jonathan Cohen. Emergent symbols through binding in
768 external memory. In *International Conference on Learning Representations*, 2020.
- 769 David Wechsler and Habuku Kodama. *Wechsler intelligence scale for children*, volume 1. Psycholog-
770 ical corporation New York, 1949.
- 771 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
772 Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models.
773 *arXiv preprint arXiv:2206.07682*, 2022.
- 774 Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms.
775 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- 776 Jingyi Xu, Tushar Vaidya, Yufei Wu, Saket Chandra, Zhangsheng Lai, and Kai Fong Ernest Chong.
777 Abstract visual reasoning: An algebraic approach for solving raven’s progressive matrices. In
778 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
779 6715–6724, 2023a.
- 780 Yudong Xu, Wenhao Li, Pashootan Vaezipoor, Scott Sanner, and Elias B Khalil. Llms and the abstrac-
781 tion and reasoning corpus: Successes, failures, and the importance of object-based representations.
782 *arXiv preprint arXiv:2305.18354*, 2023b.
- 783 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,
784 Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. *arXiv preprint*
785 *arXiv:2407.10671*, 2024.
- 786 Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Li-
787 juan Wang. The dawn of llms: Preliminary explorations with gpt-4v (ision). *arXiv preprint*
788 *arXiv:2309.17421*, 9(1):1, 2023.
- 789 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan.
790 Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural*
791 *Information Processing Systems*, 36, 2024.
- 792 Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng
793 Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. *arXiv preprint*
794 *arXiv:2403.04652*, 2024.
- 795 Mert Yuksekogul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and
796 why vision-language models behave like bags-of-words, and what to do about it? In *The Eleventh*
797 *International Conference on Learning Representations*, 2022.
- 798 Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual
799 commonsense reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and*
800 *pattern recognition*, pp. 6720–6731, 2019.
- 801 Aimen Zerroug, Mohit Vaishnav, Julien Colin, Sebastian Musslick, and Thomas Serre. A benchmark
802 for compositional visual reasoning. *Advances in neural information processing systems*, 35:
803 29776–29788, 2022.

810 Chi Zhang, Feng Gao, Baoxiong Jia, Yixin Zhu, and Song-Chun Zhu. Raven: A dataset for relational
811 and analogical visual reasoning. In *Proceedings of the IEEE/CVF conference on computer vision*
812 *and pattern recognition*, pp. 5317–5327, 2019.

813 Daoan Zhang, Junming Yang, Hanjia Lyu, Zijian Jin, Yuan Yao, Mingkai Chen, and Jiebo Luo. Cocot:
814 Contrastive chain-of-thought prompting for large multimodal models with multiple image inputs.
815 *arXiv preprint arXiv:2401.02582*, 2024a.

816 Duzhen Zhang, Yahan Yu, Chenxing Li, Jiahua Dong, Dan Su, Chenhui Chu, and Dong Yu. Mm-llms:
817 Recent advances in multimodal large language models. *arXiv preprint arXiv:2401.13601*, 2024b.

818 Yizhe Zhang, He Bai, Ruixiang Zhang, Jiatao Gu, Shuangfei Zhai, Josh Susskind, and Navdeep Jaitly.
819 How far are we from intelligent visual deductive reasoning? *arXiv preprint arXiv:2403.04732*,
820 2024c.

821 Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal
822 chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023.

823 Kai Zhao, Chang Xu, and Bailu Si. Learning visual abstract reasoning through dual-stream networks.
824 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 16979–16988,
825 2024.

826 Liang Zhou, Kevin A Smith, Joshua B Tenenbaum, and Tobias Gerstenberg. Mental jenga: A coun-
827 terfactual simulation model of causal judgments about physical support. *Journal of Experimental*
828 *Psychology: General*, 152(8):2237, 2023.

829 Qiji Zhou, Ruochen Zhou, Zike Hu, Panzhong Lu, Siyang Gao, and Yue Zhang. Image-of-thought
830 prompting for visual reasoning refinement in multimodal large language models. *arXiv preprint*
831 *arXiv:2405.13872*, 2024.

832 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: En-
833 hancing vision-language understanding with advanced large language models. *arXiv preprint*
834 *arXiv:2304.10592*, 2023.

835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

864	Appendices	
865		
866		
867	CONTENTS	
868		
869	1 Introduction	1
870		
871	2 Related Works	3
872		
873	3 MaRs-VQA Dataset	4
874		
875	4 Visual Cognition Benchmark (VCog-Bench)	4
876		
877	4.1 Multi-Image Reasoning via Chain-of-Thought (CoT)	4
878	4.2 Vision-Language Models (VLMs)	5
879		
880	5 Experiments	6
881		
882	5.1 Experimental Settings	6
883	5.2 Experimental Results	7
884	5.3 Ablation Study	8
885	5.4 Qualitative Analysis	9
886	5.5 Visualization	10
887		
888	6 Discussion	10
889		
890	7 Conclusion	10
891		
892	Appendices	17
893		
894	A Datasets & Benchmarking Code	18
895		
896	B Data Collection and Licenses	18
897		
898	C Experimental Settings	19
899		
900	C.1 Implementation Details	19
901	C.2 Difficulty Levels in MaRs-VQA	20
902	C.3 More Qualitative analysis	21
903	C.4 System Prompts	21
904		
905	D Further Discussion on Limitations and Future Work	21
906		
907	E Ethics Discussion	23
908		
909	E.1 Negative Societal Impacts	23
910	E.2 Mitigating Bias and Negative Societal Impacts	23
911		
912		
913		
914		
915		
916		
917		

A DATASETS & BENCHMARKING CODE

We release the data and annotations of MaRs-VQA anonymously:

huggingface.co/datasets/vcog/marsvqa

We also release the initial version of code for MLLM inference in an anonymous github repo:

anonymous.4open.science/r/VCog-Bench-94D2

B DATA COLLECTION AND LICENSES

We showed and compared all datasets in VCog-Bench in Table 7. The data collection of VCog-Bench follows strict procedures. The reason we choose RAVEN, CVR, MaRs-VQA is because all these datasets contain zero-shot / few-shot human investigation results. Based on these results, we can compare the MLLM’s performance with human in matrix reasoning tasks.

For RAVEN and CVR, we followed the original data generation pipeline in their repo. For MaRs-VQA, we download all questionnaires from MaRs-IB and then re-annotate all images by ourselves.

RAVEN The original dataset link of RAVEN is github.com/WellyZhang/RAVEN. It is under GPL-3.0 License (RAVEN LICENSE) and is free to use by public. All data in RAVEN are generated by rule-based scripts. We follow the basic setting of RAVEN, and modify the range of COLOR_VALUES to [255, 192, 128, 64, 0] and SIZE_VALUES to [0.3, 0.45, 0.6, 0.75, 0.9]. The sample size of RAVEN in VCog-Bench is 560.

CVR The original dataset link of CVR is github.com/serre-lab/CVR. It is under Apache License 2.0 (CVR LICENSE). CVR is an accepted paper by NeurIPS 2022 Datasets and Benchmarks track, so all of its data is free to use by public. We follow the same data generation pipeline in CVR to generate 309 samples.

MaRs-VQA The image data of MaRs-VQA is from MaRs-IB (Chierchia et al., 2019) and annotated with context option by our team. It contains 18 questionnaires, each of questionnaire contains 80 matrix reasoning questions. The human study of MaRs-IB is rigorous. In MaRs-IB’s original user study, all participants provided informed consent and all procedures were approved by UCL’s ethical committee.

The paper and study results are under MIT License. All questionnaires are under Attribution-NonCommercial 3.0 (MaRs-IB LICENSE), which means it allows people to use the work, or adaptations of the work, for noncommercial purposes only, and only as long as they give credit to the creator. Thus, the MaRs-VQA dataset will under the same license.

After we download all questionnaires from MaRs-IB, we use two Python scripts to merge all question-option pairs from different questionnaires into the same sample set. Then, we generate Option Set A, Option Set B in Figure 6 by manipulating the size and image position of option images. After that, we annotate the language description of 4 options in 10 samples from the raw data. The language description is used as system prompt to guide GPT-4o to generate all description for all data in MaRs-VQA. Then, human annotators review the annotation and revise them. Finally, we publish all annotations as Option Set A, Option Set B, and Option Set C for MaRs-VQA. Figure 6 shows an example of each type of option.

The sub-task statistics of MaRs-VQA is in Table.

Compared to other zero-shot matrix reasoning dataset (Table 1) to evaluate matrix reasoning for MLLMs, MaRs-VQA has advantages list below:

- MaRs-VQA comprises 1,440 image instances designed by psychologists, making it the largest dataset for zero-shot matrix reasoning evaluation.
- MaRs-VQA includes a diverse range of data, such as variations in color, geometry, positional relationships, and counting.

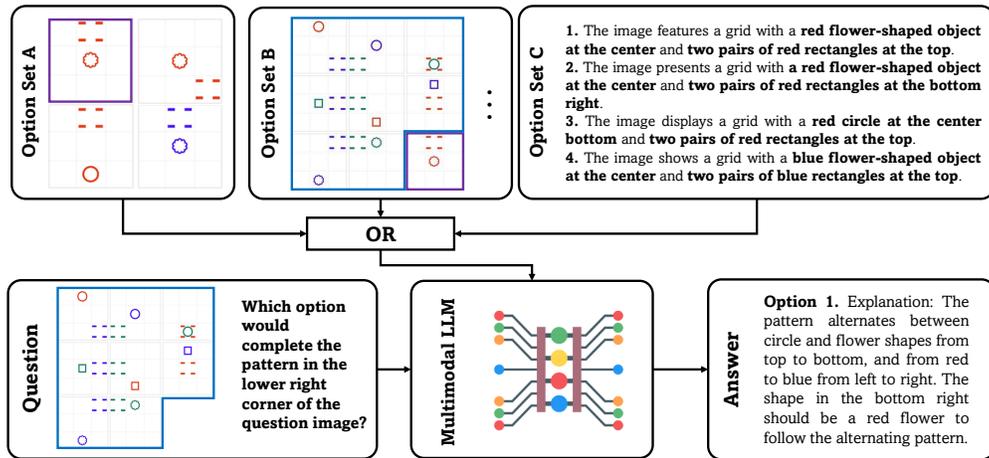


Figure 6: VQA Design of MaRs-VQA to evaluate Multimodal LLMs. The input set contains an image with a corresponding question and three sets of four-option images/contexts. Option Set A includes single-object images that can be filled into the blank region. Option Set B includes full 3x3 images containing all objects. Option C includes language descriptions for each option.

- The data source for MaRs-VQA is MaRs-IB (Chierchia et al., 2019), which is based on rigorous human studies. This dataset is widely recognized in the psychology community and has inspired numerous follow-up studies in child psychology and pediatrics. This is the first time we introduce it to the AI/ML community.

C EXPERIMENTAL SETTINGS

Dataset	Question	Option	Instance	Description
RAVEN (Zhang et al., 2019)			rule-based generation	8 options per instance grayscale image rule-based stimuli include human study
CVR (Zerroug et al., 2022)	Find the outlier among 4 images		rule-based generation	4 options per instance RGB image rule-based stimuli include human study
MaRs-VQA			1,440	4 options per instance RGB image psychologist designed stimuli include human study

Table 7: Datasets in the VCog-Bench. Both the RAVEN and CVR are rule-based generated datasets. The test samples in MaRs-VQA are designed by psychologists from MaRs-IB.

C.1 IMPLEMENTATION DETAILS

We used langchain to implement all closed-source MLLMs. The temperature of all models are 0 and the max token length is 1024. For all datasets, we follow their default image size, type settings for closed-source MLLMs. All experiments are run with three different random seeds, however, since we set temperature to 0, the final accuracy is the same for all random seeds.

For open-source models, we use the public available weights and data loader settings from the HuggingFace. InstructBLIP (Dai et al., 2024) and MiniGPT-4 (Zhu et al., 2023) are used their original GitHub repo to implement the zero-shot matrix reasoning inference pipeline. Testing is conducted using two NVIDIA RTX 4090 GPUs for 7B-sized VLMs and eight NVIDIA A100 80GB GPUs for VLMs larger than 7B. All experiments are run with three different random seeds, and the results are averaged.

C.2 DIFFICULTY LEVELS IN MARS-VQA

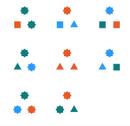
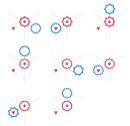
Difficulty Level	Question	Option	Description
1			Shape + Size
2			Color + Multi-object
3			Shape + Color + Position
4			Shape + Color + Multi-object
5			Shape + Color + Position + Multi-object

Table 8: Explanation of Difficulty Levels.

Based on Figure 8, here is the explanation of difficulty levels presented in our paper:

- **Difficulty Level 1:** Single sub-task and two simple sub-tasks Description: The task involves only one changing attribute across the matrix reasoning—either shape, color, size, position, or multi-object. Or two simple attributes: (shape & color), (shape & size), (shape & position), (color & size), (color & position), (size & position). Example: Figure 4 (top-left) is a matrix reasoning task where only the size and color of the objects changes. This is a difficulty level 1 task.
- **Difficulty Level 2:** Two sub-tasks involving multi-object sub-task Description: The task involves multiple objects combined with one other changing attribute. The sub-task combinations are (multi-object & shape), (multi-object & color), (multi-object & size), (multi-object & position).
- **Difficulty Level 3:** Three simple sub-tasks combined Description: The task involves three changing attributes simultaneously. The sub-task combinations are (shape & color & size), (shape & position & size), (shape & position & color), (size & position & color).
- **Difficulty Level 4:** Three sub-tasks involving multi-object sub-task Description: The task involves multiple objects combined with two other changing attributes. The sub-task combinations are (multi-object & shape & color), (multi-object & shape & size), (multi-object & shape & position), (multi-object & color & position), (multi-object & color & size), (multi-object & position & size).

- **Difficulty Level 5 and Above:** Four or more Sub-tasks Description: The task involves combinations of four or five attributes. Example: Figure 4 (top-right) is a matrix reasoning task (shape & position & color & multi-objects) and its difficulty level is > 4 .

As more attributes change simultaneously, the task becomes more complex, requiring higher levels of abstract reasoning to identify patterns. In addition, each additional changing element adds to the cognitive load, making it more challenging to discern the correct answer.

C.3 MORE QUALITATIVE ANALYSIS

In this section, we further analyze the failure cases of GPT-4o. Correct reasoning is highlighted in green, while incorrect reasoning is marked in red. Although GPT-4o is sometimes able to extract a subset of key information from the question image, it frequently fails to arrive at the correct final answer. This is primarily due to critical features being either overlooked or inadequately utilized in the decision-making process. As a result, the final answers are often incorrect or only partially aligned with the relevant attributes. It reveals that visual working memory will be a key part to optimize the MLLM’s performance in matrix reasoning problem.

C.4 SYSTEM PROMPTS

For each dataset, we prepare custom system prompt. Their pipeline is similar. First, we created a system message prompt (see Figure 8, 9 for zero-shot inference, and Figure 10, 11, 12 for CoT) to guide the MLLM understanding the basic information of matrix reasoning tasks and the structure of the input, and formulating multiple-option images or contexts. The difference for zero-shot and CoT is we provide the guideline to encourage the model think the problem step-by-step based on extracting all useful information from structure $K = \{[r, a, o, s] | r \in \mathcal{R}, a \in \mathcal{A}, o \in \mathcal{O}, s \in \mathcal{S}\}$. The output format is a json structure including “Answer” and “Explanation” as keys.

D FURTHER DISCUSSION ON LIMITATIONS AND FUTURE WORK

Insights Unlike other VQA benchmarks, our work approaches the perspective of human visual cognition—an underexplored domain. Based on our experimental results, we offer the following insights for vision researchers:

- While scaling laws have some applicability to visual cognition tasks, merely increasing model size and training data is insufficient to achieve human-level performance.
- To demonstrate that VLMs possess strong visual cognitive abilities, it is crucial to evaluate them on zero-shot inference tasks like matrix reasoning—tasks characterized by simple visual content but requiring complex reasoning to find the correct answer.
- Unlike other multi-image visual reasoning benchmarks, VCog-Bench effectively highlights the performance gap between MLLMs and human cognition in these tasks.

From our main and ablation experiments, we observed that as task difficulty increases, the performance of MLLMs in multi-image reasoning scenarios deteriorates. Interestingly, providing language-based descriptions of each option (i.e., inputting the model with a single question image and context-based options) improved the models’ performance compared to using multi-image options. This suggests that language still plays a significant role in the visual reasoning processes of current MLLMs and VLMs.

In contrast, human visual cognition—especially in children—allows individuals to solve matrix reasoning tasks without relying on advanced language reasoning capabilities. Children can often solve these tasks effectively by utilizing their visual working memory and pattern recognition skills.

One potential reason for the performance gap is that current MLLMs/VLMs may underemphasize the visual encoder relative to the language encoder. In many recently released VLMs, the visual module is much smaller than the language model module, and the visual encoders are frozen during Large Language Model (LLM) and alignment layer fine-tuning in open-sourced VLMs. This imbalance might limit the models’ capacity to retain and process complex visual information during reasoning tasks.

1134 To better retain visual information during the reasoning process, MLLMs may require more capable
1135 visual modules that can handle complex visual patterns and maintain this information throughout
1136 the reasoning steps. Moreover, optimizing the training process with end-to-end multimodal training—
1137 without freezing any layers in the visual modules—can be beneficial. Recent models have
1138 begun to explore end-to-end VLM fine-tuning, demonstrating the potential of this approach, though
1139 challenges remain such as the need for multi-round alignment. In the future, developing more
1140 advanced methods to effectively integrate visual and linguistic features will be crucial.

1141
1142 **Limitations** In the main paper, we briefly discussed the limitations of our work. Here, we provide
1143 a more in-depth discussion. First, our dataset is composed of limited publicly available matrix
1144 reasoning datasets, which must include human study results. The RAVEN and CVR datasets, created
1145 by the AI/ML community, were not developed following rigorous psychological research norms.
1146 Consequently, our benchmarking results, which utilize these datasets, should not be used to derive
1147 psychological or clinical conclusions. While MaRs-VQA addresses this problem, its samples cannot
1148 represent all formats of matrix reasoning found in IQ tests such as the WISC and the Cattell Culture
1149 Fair Intelligence Test (Cattell & Cattell, 1960). We cannot use these IQ tests directly because they
1150 are not freely available, and copyright restrictions usually prevent these pen-and-paper tasks from
1151 being adapted into computerized formats.

1152 Second, the size of the datasets in VCog-Bench is relatively small compared with typical computer
1153 vision datasets, due to the inherent challenges involved in collecting matrix reasoning data. However,
1154 as we have argued in our paper, matrix reasoning should not be presented in typical machine learning
1155 settings—fine-tuning models on training sets and evaluating performance on test sets. Benchmarking
1156 MLLMs’ visual reasoning performance should be conducted in a zero-shot inference setting, ensuring
1157 that all data in the test set are not included in the models’ training data. Even compared with other
1158 recently released human-designed matrix reasoning datasets, ours is still the largest (see Table 1).

1159
1160 **Future Work** Although LLMs have achieved remarkable success in language understanding and
1161 generation, a significant portion of their parameters is dedicated to encoding linguistic patterns and
1162 memorizing factual information, which offers limited benefits for tasks requiring visual cognition.
1163 This disparity between Multimodal LLMs and humans indicates that merely increasing model size is
1164 insufficient to achieve human-level zero-shot inference in these domains. While our benchmark and
1165 baseline models represent a significant initial step, further data collection and in-depth human studies
1166 remain essential.

1167 From our experimental results, we observe that current MLLMs have enhanced basic matrix reasoning
1168 capabilities, with models like GPT-4o and Gemini Pro 1.5 achieving significantly higher accuracy
1169 than random guessing across all three matrix reasoning tasks. By using Monte Carlo Tree Search
1170 to optimize the results via multi-round reasoning and exclusion, GPT-4o can achieve much better
1171 outcomes, albeit at the cost of increased inference time. We anticipate that the next generation
1172 of MLLMs will approach human-level performance in matrix reasoning. It is crucial to maintain
1173 these visual cognition-based benchmarks, continuously monitor the performance of newly released
1174 MLLMs, and encourage open-source MLLMs and VLMs to include matrix reasoning tasks for
1175 performance comparison.

1176 Finally, we pose the open-ended question of whether MLLMs need to achieve or surpass human-level
1177 zero-shot inference capability in matrix reasoning tasks. Addressing this issue requires drawing
1178 on theories from cognitive science and psychology to understand the nature of human and MLLM
1179 intelligence. Matrix reasoning ability develops early in human neurodevelopment, with children as
1180 young as four providing sensible answers to simple matrix reasoning questions without additional
1181 training, making it a critical component of IQ tests. In contrast, LLMs and MLLMs rely on training
1182 data, fundamentally differing from how children develop cognitive abilities. However, we believe that
1183 these two learning processes share commonalities: both involve the gradual accumulation of skills
1184 and the ability to generalize from past experiences. Exploring these parallels can provide valuable
1185 insights into designing MLLMs that more closely mimic human visual cognition, ultimately leading
1186 to more advanced and capable models. Additionally, we observe that current open-source models
1187 achieve matrix reasoning performance very close to that of closed-source models. However, VLMs
face challenges in supporting multiple images as input and managing visual memory. Addressing
these challenges is a crucial direction for building more robust open-source VLMs in the future.

1188 E ETHICS DISCUSSION

1189

1190 This research aims to advance LLMs and VLMs by providing a new benchmark for evaluating AI
1191 capabilities in visual reasoning. MaRs-VQA builds on the MaRs-IB (Attribution-NonCommercial
1192 3.0 License), and VCog-Bench builds on MaRs-VQA, RAVEN (GPL-3.0 License), CVR (Apache
1193 License 2.0). All code and data are available on GitHub. No conflicts of interest exist among the
1194 study’s contributors. More discussion on the ethical aspects of VCog-Bench is included in the
1195 Appendix. The annotation process is IRB approved by a clinical institute.

1196

1197 E.1 NEGATIVE SOCIETAL IMPACTS

1198

1199 We foresee no direct negative societal impacts from our matrix reasoning benchmark. However, it
1200 could be misunderstood or misinterpreted as comparing AI “thought” to human cognition or misused
1201 to evaluate human abilities across demographics or ethnicity. We strongly caution against such
1202 misuse, as our datasets are not validated for human assessment.

1203 Another concern relates to the future conclusion from our benchmark. While matrix reasoning is
1204 a crucial test for evaluating human intelligence, observing that VLMs with large model weights
1205 perform better on matrix reasoning tasks does not imply that the intelligence of MLLMs follows the
1206 same “scaling law” from the general domain. A comprehensive intelligence test requires accurate
1207 assessment using human-based tools, of which matrix reasoning is only one critical component. We
1208 cannot conclude that larger MLLMs can achieve human intelligence.

1209 Additionally, there is a potential concern for discrimination against certain groups based on race,
1210 gender, or age in human study results. Although all human results in our experiment tables are
1211 sourced from previously published papers, we cannot guarantee that all previous research adhered to
1212 strict standards ensuring the inclusion of all groups in the human investigation process.

1213

1214 E.2 MITIGATING BIAS AND NEGATIVE SOCIETAL IMPACTS

1215 While the use of VCog-Bench and MaRs-VQA come with potential negative social impacts, there
1216 are viable mitigations that can address these concerns. These include adding instructions for proper
1217 use and restricting unethical human investigations. Users must be aware of the ethical implications
1218 associated with our benchmark and take appropriate measures to ensure its safe and responsible
1219 utilization.

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Question Image	Options	Reasoning
		In each row, the top shape is the same, and the bottom two shapes are always pink. Following this pattern, the missing shape should have the same top shape as the others in its row, which is pink, with two pink shapes at the bottom.
		Each row contains the same three shapes in different positions. The missing image should include the red circle with a blue plus, a blue circle, and a red triangle. Option 3 matches this pattern.
		Each row has a pattern where the number of blue dots increases by one as you move from left to right. The first column has 4 dots, the second has 5, and the third has 6. Following this pattern, the missing shape should have 6 dots.
		In each row, the number of dots in the top part increases from left to right (1, 2, 3). Following this pattern, the bottom-right cell should have three dots on the top, matching the third column's pattern for the row.
		In each row, a circle and a pair of parallel lines appear, with the circle shifting positions and colors. The third column should have a pink circle and pink lines to complete the pattern in the last row.
		Each row follows a pattern where the main shapes alternate positions, with the half-circle and ellipses appearing in different orientations. The third row should mirror this pattern. In the first row, the half-circle is on the left, in the second row, it is on the right, and in the third row, it should be in the middle.

Figure 7: More visualization results for GPT-4o's reasoning.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

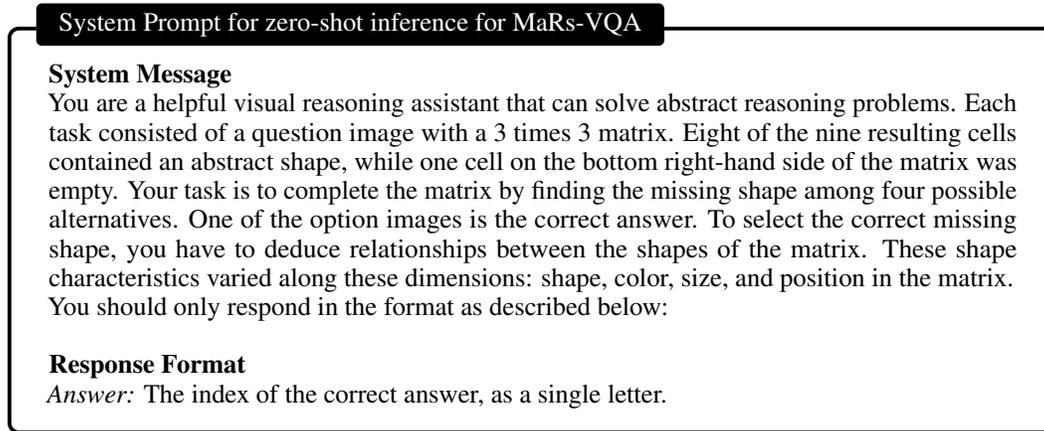


Figure 8: System prompts for zero-shot MLLM inference of MaRs-VQA.

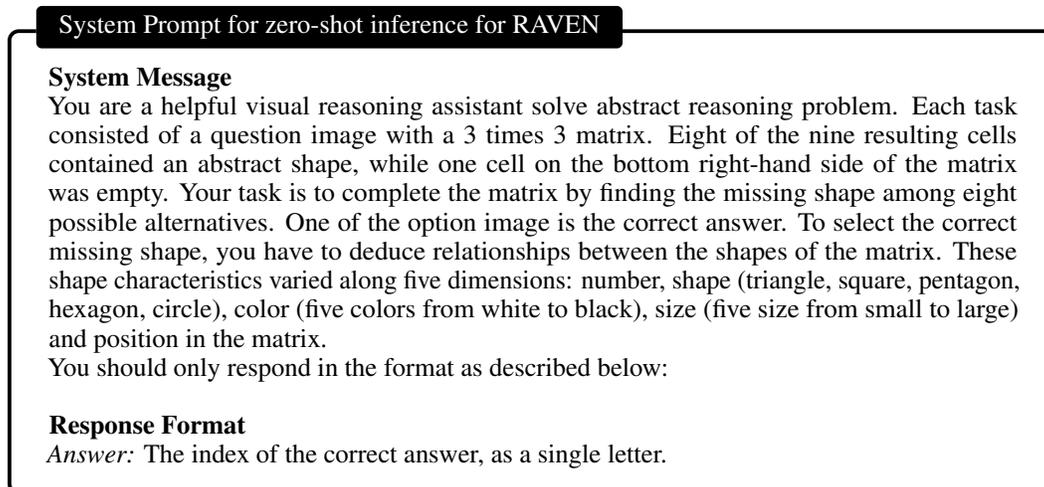


Figure 9: System prompts for zero-shot MLLM inference of RAVEN.

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

System Prompt for MLLMs with CoT for MaRs-VQA

System Message

You are a helpful visual reasoning assistant that can solve abstract visual reasoning problems. Each task consisted of a question image with a 3 times 3 matrix. Eight of the nine resulting cells contained an abstract shape, while one cell on the bottom right-hand side of the matrix was empty. Your task is to complete the matrix by finding the missing shape among four possible alternatives. One of the options is the correct answer.

The first step is to describe what is the attribute and relationship between each attribute in each cell of the 3 times 3 question image. The attributes can be number, position, shape, size, and color. The cell may contain multiple attributes. The relation might be '3 times 3 sub-blocks', 'rotation', 'insideness'.

The second step is to summarize the relation of three patterns in the first row of the question image, the relation of three patterns in the second row of the question image, the relation of two patterns in the third row of the question image.

Answer this question: What are the row-based high-order rules in the question image?

Based on the description for each option, answer this question: What is the constraint of all options?

Finally, infer what are the potential attributes, objects, relations in the missing cell?

You should only respond in the format as described below:

Response Format

Explanation: The step-by-step reasoning for the answer.

Answer: The index of the correct answer, as a single letter.

Figure 10: System prompts for MLLM CoT inference of MaRs-VQA.

System Prompt for MLLMs with CoT for RAVEN

System Message

You are a helpful visual reasoning assistant that can solve abstract visual reasoning problems. Each task consisted of a question image with a 3 times 3 matrix. Eight of the nine resulting cells contained an abstract shape, while one cell on the bottom right-hand side of the matrix was empty. Your task is to complete the matrix by finding the missing shape among eight possible alternatives. One of the option images is the correct answer.

The first step is to summarize the relation of three patterns in the first row of the question image, the relation of three patterns in the second row of the question image, the relation of two patterns in the third row of the question image. What is this relation? The features in the patterns can be constant, progression, arithmetic, distribute three. Try to describe this relationship.

The second step is to describe what is the attribute and relationship between each attribute in each cell of the 3 times 3 cells question image and four option images. The attributes can be number; shape (triangle, square, pentagon, hexagon, circle); colour (five colors: white, light gray, gray, dark gray, black); size (five size: tiny, small, medium, large, huge); and positional relation (inside outside relation, left right relation, top down relation, two times two sub-blocks, 3 times 3 sub-blocks). The cell may contain multiple attributes.

Finally, give me the answer based on step 1-2.

You should only respond in the format as described below:

Response Format

Explanation: The step-by-step reasoning for the answer.

Answer: The index of the correct answer, as a single letter.

Figure 11: System prompts for CoT MLLM inference of RAVEN.

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

System Prompt for MLLMs with CoT for CVR

System Message

You are a helpful visual reasoning assistant that can solve abstract reasoning problems. Each task consisted of four option images. Your task is to identify which image is different from the other three? To find the correct answer, you have to deduce relationships inside each image and then find the difference.

The first step is to describe what is the attribute and relationship between each attribute in four option images. The features can be number of objects; shape; color; size, relationship of colour or shape or size or direction among objects; and positional relation of objects (inside outside relation, left right relation, top bottom relation, adjacent relation). Each image may contain multiple attributes and multiple relations.

Based on the description and image for each option, answer this question: What is the constraint / similarity of most of the options?

Finally, infer which image is the outlier?

You should only respond in the format as described below:

Response Format

Explanation: The step-by-step reasoning for the answer.

Answer: The index of the correct answer, as a single letter.

Figure 12: System prompts for CoT MLLM inference of CVR.