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 closedsource 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

030 031

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029

032 Matrix reasoning is a crucial ability in human perception and cognition, essential for nonverbal, 033 culture-reduced intelligence measurements as it can minimize the influence of acquired knowledge 034 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 & 035 Mańdziuk, 2023) (see Figure 1). Participants have to identify and comprehend the rules based on a few 036 provided patterns, and then reason about the next pattern following the same rules. Matrix reasoning 037 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 040 methods have been proposed as a part of fluid intelligence tests. The two most famous assessments 041 are Wechsler Intelligence Scale for Children (WISC) (Wechsler & Kodama, 1949) and Raven's 042 Progressive Matrices (RPM) (Raven, 2003). 043

Recently, matrix reasoning tasks have emerged as an ideal testbed for investigating whether deep 044 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., 046 2018; Zhang et al., 2019; Webb et al., 2020). Previous research on matrix reasoning assessments 047 applied typical machine learning settings – finetuning models on training sets and evaluating the 048 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 051 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 052 everyday experiences in natural scenes. Previous machine learning models ignore these prerequisites when modeling matrix reasoning problem. This could lead to an overestimation of the models'



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.

reasoning abilities, as they might be leveraging learned patterns specific to the training data rather
 than demonstrating genuine generalization and reasoning skills from visual cognition.

077 Recently, Multimodal Large Language Models (MLLMs) have shown surprising understanding and reasoning capabilities, marking an important milestone towards Artificial General Intelligence 079 (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 081 exposure by in-context learning. However, current MLLMs remain inadequate in visual cognition 082 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 084 on abstract reasoning skills. The RAVEN IQ-test also has some limitations, including a small 085 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. 087

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 090 psychologist-designed dataset for matrix reasoning assessment including 1,440 examples in total. 091 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 092 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) 094 and human performance under a zero-shot inference setting (no prior knowledge) on MaRs-VQA 095 and other abstract reasoning datasets containing human studies. In our experiments, we observe that 096 MLLMs with more parameters generally perform better on our benchmark, adhering to established 097 scaling laws in a limited scope. However, even the largest open-source MLLMs and GPT-40 fall 098 short of surpassing human performance in matrix reasoning tasks. Furthermore, many MLLMs have 099 a mismatch in performance between matrix reasoning tasks and other general VQA benchmarks, 100 which provides some insights into the drawbacks of existing models. In conclusion, our contributions 101 are summarized as follows:

- 102 103
- 104 105

• 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	×	×	1	×	×
Visual Reasoning Benchmark (COLM-24) (Zhang et al., 2024c)	Mensa Test, RAVEN, IntelligenceTest		241	×	×	×	×	×
MaRs-VQA (ours)	MaRs-IB	• · · · · · · · · · · · · · · · · · · ·	1,440	1	1	1	1	1

126 127 128

129

125

108

109

110

111 112

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 130 exploring human-like AI in psychology and cognition (Ullman, 2023). Recent works tested LLMs' 131 cognitive abilities in causal reasoning (Binz & Schulz, 2023), abstract reasoning (Xu et al., 2023b; 132 Moskvichev et al., 2023; Jiang et al., 2024b; Ahrabian et al., 2024), analogical reasoning (Webb et al., 133 2023), systematic reasoning (Hagendorff et al., 2023), and theory of mind (Strachan et al., 2024). 134 Their observation showed that LLMs like GPT-4 (Achiam et al., 2023) have been proven successful in 135 most cognitive tests related to language-based reasoning. Despite this success, only limited research 136 has been conducted on the areas of MLLMs and visual cognition. Visual cognition involves the 137 process by which the human visual system interprets and makes inferences about a visual scene using 138 partial information. Buschoff et al. observed that while LLMs demonstrate a basic understanding of 139 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 140 abilities in simple tasks (Lerer et al., 2016; Zhou et al., 2023; Jassim et al., 2023), and ignore testing 141 complex abstract reasoning and logical reasoning ability related to fluid intelligence. Therefore, new 142 and challenging benchmarks based on the theory of visual cognition are needed to assess and improve 143 AI systems' capabilities for human-like visual understanding. 144

145 Matrix Reasoning Matrix reasoning is often used to determine human intelligence related to visual 146 cognition and working memory (Salthouse, 1993; Jaeggi et al., 2010; Fleuret et al., 2011) that is 147 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 148 among visual objects and use reasoning to find visual cues. Early research indicated that deep 149 learning models can be trained with large-scale matrix reasoning datasets to solve simple matrix 150 reasoning (Stabinger et al., 2021; Małkiński & Mańdziuk, 2022; 2023; Xu et al., 2023a; Małkiński & 151 Mańdziuk, 2024) and compositional visual relation tasks (Fleuret et al., 2011; Zerroug et al., 2022; 152 Ommer & Buhmann, 2007; Liu et al., 2021), achieving human-level accuracy. Several datasets and 153 benchmarks are also proposed, such as PGM (Barrett et al., 2018), RAVEN (Zhang et al., 2019), 154 RAVEN-I (Hu et al., 2021), RAVEN-FAIR (Benny et al., 2021), CVR (Zerroug et al., 2022). However, 155 these works have a key limitation. They ignore that humans can solve these problems by zero-shot 156 reasoning without explicitly learning from large-scale data. After the blooming of LLMs, researchers 157 are keen on testing whether LLMs reached the same abstract reasoning capabilities as humans. Webb 158 et al. (Webb et al., 2023) encode matrix reasoning into a symbolic problem based on human's prior 159 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 160 community, such as RAVEN-IQ (Huang et al., 2024) containing 50 instances, Visual Reasoning 161 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-165 Language Models (VLMs) for addressing vision reasoning tasks (Zellers et al., 2019; Bordes et al., 166 2024). These latest VLMs are constructed using a combination of the CLIP vision encoder, pre-167 trained LLMs, and a connected adapter to align visual features with language space (Zhang et al., 168 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), 170 CogVLM (Wang et al., 2023) underscore the significance of employing high-quality visual instruction 171 tuning data. Additionally, tool learning methods have also explored the potential of integrating code 172 generation pipelines with visual inference (Surís et al., 2023). Nevertheless, current VLMs encounter 173 challenges in adapting to high-resolution and visually complex images. These problems stem from 174 the absence of a robust visual search mechanism (Wu & Xie, 2023), few-shot reasoning (Guo et al., 175 2023), compositional understanding (Yuksekgonul et al., 2022) and the constrained visual grounding 176 capabilities inherent in CLIP (Tong et al., 2024).

177 178

3 MARS-VQA DATASET

179 180 181

The MaRs-VQA dataset is designed to evaluate the zero-shot abstract reasoning capabilities of
 MLLMs through various matrix reasoning VQA tasks. All sample images in MaRs-VQA are sourced
 from the questionaires from 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. We create annotations in the images from all questionaires and design the VQA annotations.

To transform the matrix reasoning problem into a VQA task, we firstly define three different option 188 sets - two image-based sets (A and B) and one language-based set (C). In Option Set A, we provide 189 four candidates to the missing patch in the question. In Option Set B, the options are created by 190 filling the four patches in Set A into the 3×3 question image. Note that Option Set B is used for 191 visualization purposes only and is not included in our experiment. We further diversify the modalities 192 of our dataset to support the evaluation of different kinds of models. Specifically, we use GPT-40 193 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 194 initial human annotations in our data collection. These examples are then used as few-shot samples to 195 query GPT-40 through in-context learning. The context generation system prompt guides GPT-40 to 196 compare all four option images and generate distinct descriptions for each one. After generating all 197 samples, human annotators in the author team review each option and revise the incorrect description. 198 Examples are showed in Figure 6 in the Appendix. 199

200 201

202 203

204

205

206

207

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.

- 208 209
- 210 211

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., 2024a;b; Mondal et al., 2024; Xu et al., 2023b) and object-centric representation learning (Seitzer

233

234

266

216 G 217 Summarize attributes, objects, relations in each Row-based high-order rules 218 row What is the aint of al 219 options Summarize attributes objects of each option A 220 Question image 222 Input Images InternVl System prompt Language options Qwen-VL **-**Visual Encode 224 VLM Word Embedding CogVLM 225 226 в D Α С Visual Language Decode 227 that can solve Each task con 228 with a 3 tim es 3 matrix. Eight of the n 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 multiple possible alternatives. One of the option images is the correct answer. E 229 Hire participants to take matrix reasoning exam. Human All participants should not see any question-answer pairs before the exam. 230 231

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.

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.

245 Based on structure K, we use three stages to guide MLLM to use human-level 246 thought to understand matrix reasoning tasks. The first stage is to guide the 247 Multimodal LLM to summarize the visual feature (e.g. shape) of each row in 248 the 3×3 question image. Then, based on these row-based visual features, the 249 model will then conclude the high-order rule/pattern \mathcal{R} . The second stage is to 250 extract the basic attributes \mathcal{A} and inner relations \mathcal{O} to integrate objects in each 251 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 253 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 254 depiction of how to leverage CoT in matrix reasoning tasks. 255

256 To further enhance CoT with diverse prompts, we introduce a multi-round 257 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 258 multi-image CoT solve the matrix reasoning problem. The selected image is 259 then incorporated into the question image as a new input, which is fed back 260 into the MLLM with a prompt directing it to evaluate the correctness of the 261 complete 3×3 matrix, specifically focusing on the bottom-right corner. If the 262 MLLM determines the result is correct, the final answer is output; otherwise, 263 the incorrect option is excluded and CoT process is repeated. 264



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). $\mathbf{x}_i = [x_i^1, ..., x_i^k]$ 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:

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

 \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 $.F_{\theta}$ is an autoregressive decoder in the LLM for answer generation. It is defined as:

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

287 In VLMs, the input question image is first processed by the visual encoder such as CLIP (Radford 288 et al., 2021). Then, additional adapter layers are used to map visual features into language feature 289 space. These features, along with the context-based option descriptions, are sent to the LLM decoder. 290 The LLM decoder then integrates the information from both the input question image and the option 291 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 292 provides a structured way to break down the problem, potentially improving interpretability compared 293 to end-to-end close-source models. 294

295 296

297

299

275 276 277

278

279

5 EXPERIMENTS

298 5.1 EXPERIMENTAL SETTINGS

Datasets In addition to MaRs-VQA, we selected two well-known open-source datasets for matrix 300 reasoning and abstract visual reasoning to conduct experiments in VCog-Bench. The first dataset is 301 RAVEN (Zhang et al., 2019), designed to probe abstract reasoning in a format similar to the Raven's 302 Progressive Matrices IQ test, with each question providing eight options. The second dataset is 303 Compositional Visual Reasoning (CVR) (Zerroug et al., 2022), which evaluates deep learning models 304 using 103 unique configurations generated by predefined rules. Each sample in CVR is an outlier 305 detection problem, with four options provided per question. However, both RAVEN and CVR share a 306 significant limitation: all samples are algorithmically generated using fixed rules, which limits their 307 diversity and lacks psychological validity.

Baselines for Multi-image Reasoning We selected the Claude 3 family (Haiku, Sonnet, Opus) (An-thropic, 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
 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

336

337

352

Method	Learning	Accuracy (%) ↑			
	Dearing	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-40 (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 (%) ↑	
				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)	unknown	unknown	unknown	34.82	35.36
GPT-40 (OpenAI, 2024)	unknown	unknown	unknown	37.38	38.84
Gemini Pro 1.5 (Reid et al., 2024)	unknown	unknown	unknown	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.

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 357 matrix reasoning problem setting, which serve as the system prompt for zero-shot inference. For 358 object-centric CoT reasoning, we create specific prompts to guide the model's thought process 359 through multiple stages, enabling step-by-step reasoning. For open-source baseline models, we use 360 the same system prompt settings across all models. Testing is conducted using two NVIDIA RTX 361 4090 GPUs for 7B-sized VLMs and eight NVIDIA A100 80GB GPUs for VLMs larger than 7B. All 362 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 (Acc = Correct/Total), 364 consistent with other VQA benchmarks (Lu et al., 2022; Liu et al., 2023).

366 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 the
Claude 3 family, GPT-4V, and GPT-40 for this task, as these models support 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,

Method	Multi-Image	Accuracy (%) ↑				
		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-40 (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
iemini Pro 1.5 (Reid et al., 2024)	×	51.10	30.21	26.19	29.17	35.51
GPT-40 (OpenAI, 2024)	×	58.89	45.83	32.14	26.39	26.09

387

388

> 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.

389 it does not benefit smaller models such as Claude 3 Sonnet. Compared to the results in MaRs-VQA and RAVEN, GPT-40 achieves much better zero-shot and object-centric CoT inference results in 390 the CVR dataset, almost matching the performance (ResNet-50: 57.9%, ViT-small: 32.7%, WReN: 391 42.4%) of fine-tuned models with 1,000 training samples in CVR's paper (Zerroug et al., 2022). 392

In the second part of our experiment, we investigated the use of VLMs (question image + language 394 options) to solve matrix reasoning problems in MaRs-VQA and RAVEN. The CVR dataset was 395 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-40 396 in MaRs-VQA and RAVEN. Notably, Gemini Pro 1.5 outperformed GPT-40 on the RAVEN dataset. 397

398 We identified three major issues after reviewing the reasoning outputs of current MLLMs in Table 2 399 and 3: (1) Limited Use of Visual Information: MLLMs cannot directly use visual features for 400 reasoning, making them insensitive to non-verbal spatial features during CoT reasoning. This 401 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 402 alone. (2) Restricted Visual Working Memory: The visual working memory of MLLMs is limited, 403 causing visual feature information to be easily lost during the text generation reasoning process. 404 (3) Integration Challenges: Even if MLLMs possess strong task-specific skills like recognition, 405 segmentation, and object detection, they struggle to integrate these skills into high-level visual 406 reasoning tasks. 407

408 409

5.3 ABLATION STUDY

410 In this subsection, we conduct ablation experiments to analyze how to improve the performance of 411 MLLMs on the matrix reasoning problem. Table 5 compares the Chain-of-Thought (CoT) baseline 412 with two approaches: few-shot reasoning and multi-round reasoning. Few-shot reasoning involves 413 providing a small number of question-answer examples alongside the CoT system prompt. Multi-414 round reasoning employs the advanced CoT strategy illustrated in Figure 3. The results show that 415 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 416 further improvement. These findings suggest that while few-shot in-context learning helps the model 417 better understand the matrix reasoning problem, it does not significantly enhance the MLLM's 418 visual reasoning capabilities for these tasks. Additionally, using a multi-round tree search improves 419 accuracy from approximately 34% to 42%, but it is considerably slower than single-round CoT, with 420 each inference taking over 30 seconds in multi-round mode. We also compare different MLLMs 421 across difficulty levels and different visual complexities in the MaRs-VQA dataset (see Table 4 and 422 Table 6). The difficulty level in our tasks is defined by the number of sub-tasks (visual complexities) 423 involved—specifically, variations in color, size, geometry (shape), positional relationships, and the 424 presence of multiple objects. The results indicate that GPT-40 exhibits difficulty sensitivity similar to 425 that of humans, whereas Claude 3 Opus does not demonstrate this ability. This is because GPT-40 426 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. 427

428 429

430

QUALITATIVE ANALYSIS 5.4

In this subsection, we use case studies from the MaRs-VQA dataset to illustrate how MLLMs fail in 431 some tasks and provide insights on how to improve MLLMs and VLMs for this task.

Strategy	Accuracy $(\%) \uparrow$
СоТ	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

460

461

462

463

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 5: Ablation on prompt selection.

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 \geq 4. The reasoning is a short summary of the CoT output, not the full version

464 First, we present an example to explain why the Claude 3 family performs worse than GPT-40 and 465 even worse than random guessing in all of our experiments. Figure 4 top left is one of the most 466 simple cases in MaRs-VQA's level 1 difficulty, Claude 3 Opus incorrectly identifies the shape as the 467 main target of this matrix, while the actual target is the size. In contrast, GPT-40 correctly discerns 468 the relationship between rows, noting: "The pattern in each row ends with a smaller shape colored 469 differently from the first two shapes." This example highlights a critical shortcoming in Claude 3 470 Opus's reasoning ability: limited Use of Visual Information, demonstrating its struggle to accurately interpret the key attributes in matrix reasoning tasks. GPT-40, on the other hand, showcases a 471 superior understanding of the relationships and patterns within simple data, leading to more accurate 472 responses. 473

474 However, the difficulty of the tasks increases, the performance of MLLMs deteriorates in multi-image 475 reasoning. Figure 4 bottom left and shows an example, it is the level 6 difficulty containing shape, 476 positional relation, shape with different objects. For these questions containing complex visual 477 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 478 utilized in selecting the correct option. Consequently, the final answers are often incorrect or only 479 partially related to the relevant attributes. It suggests that MLLMs are affected by the cognitive load 480 associated with processing multiple sub-tasks simultaneously, which is closely related to the concept 481 of visual working memory. The right two examples of Figure 4 also present the same observation. 482 Additionally, we observed that GPT-40 is not sensitive to the positional relationships for multi-objects 483 in the question images. 484

485 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

490

491

492

493

494

495

496

497

498 499

500

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

501 Social Impacts In the present work, we em-502 phasize that zero-shot matrix reasoning is a key item to validate human-level intelligence, 504 though it is still unclear how matrix reasoning 505 ability is acquired early in human neurodevelopment. Children's visual reasoners (without 506 any additional training) can provide sensible an-507 swers to matrix reasoning questions as early as 508 age four. The long-term goal of our work is 509 twofold. The first one is to explore the problem 510 of how close AIs or MLLMs are to human-like



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).

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-518 level zero-shot inference capability in matrix reasoning tasks. Addressing this issue requires new 519 theories from cognitive science and psychology to accurately evaluate and compare human and 520 MLLM intelligence. Unlike MLLMs, which rely on training data and domain-specific skills, human 521 cognition develops gradually and evolves with age. Humans can also learn how to solve the problem 522 progressively from previous seen matrix reasoning tasks while they are taking the test, but MLLM 523 can not learn from it via in-context learning due to the maximum tokens length. Therefore, AI 524 researchers, psychologists, and cognitive scientists must collaborate to rethink how to benchmark 525 MLLM intelligence with human intelligence.

526 527

517

7 CONCLUSION

528 529

We introduce VCog-Bench, a publicly available zero-shot matrix reasoning benchmark designed 530 to evaluate the visual cognition capability and intelligence of Multimodal Large Language Models 531 (MLLMs). This benchmark integrates two well-known datasets RAVEN and CVR from the AI 532 community and includes our newly proposed MaRs-VQA dataset. We also introduce several important 533 concepts to redefine zero-shot matrix reasoning task evaluation, focusing on multi-image reasoning 534 with object-centric Chain-of-Thought (CoT) system prompts. Our findings show that current state-ofthe-art MLLMs and Vision-Language Models (VLMs), such as GPT-40 and LLaVA-1.6, InternVL 536 demonstrate some basic understanding of matrix reasoning tasks. However, these models still face 537 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 538 aim to encourage further innovation and progress in the field of improving the visual cognition of MLLMs.

540 REFERENCES

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.

- Kian Ahrabian, Zhivar Sourati, Kexuan Sun, Jiarui Zhang, Yifan Jiang, Fred Morstatter, and Jay
 Pujara. The curious case of nonverbal abstract reasoning with multi-modal large language models. *arXiv preprint arXiv:2401.12117*, 2024.
- 548 549 550 Anthropic. Introducing the next generation of claude. https://www.anthropic.com/news/ claude-3-family, 2024.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang
 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities.
 arXiv preprint arXiv:2308.12966, 2023.
- David Barrett, Felix Hill, Adam Santoro, Ari Morcos, and Timothy Lillicrap. Measuring abstract reasoning in neural networks. In *International conference on machine learning*, pp. 511–520.
 PMLR, 2018.
- Yaniv Benny, Niv Pekar, and Lior Wolf. Scale-localized abstract reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12557–12565, 2021.
- Marcel Binz and Eric Schulz. Using cognitive psychology to understand gpt-3. *Proceedings of the National Academy of Sciences*, 120(6):e2218523120, 2023.
- Florian Bordes, Richard Yuanzhe Pang, Anurag Ajay, Alexander C. Li, Adrien Bardes, Suzanne Petryk, Oscar Mañas, Zhiqiu Lin, Anas Mahmoud, Bargav Jayaraman, et al. An introduction to vision-language modeling. *arXiv preprint arXiv:2405.17247*, 2024.
- Zheng Cai, Maosong Cao, Haojiong Chen, Kai Chen, Keyu Chen, Xin Chen, Xun Chen, Zehui Chen,
 Zhi Chen, Pei Chu, et al. Internlm2 technical report. *arXiv preprint arXiv:2403.17297*, 2024.
- Giacomo Camposampiero, Loïc Houmard, Benjamin Estermann, Joël Mathys, and Roger Wattenhofer.
 Abstract visual reasoning enabled by language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2642–2646, 2023.
- Patricia A Carpenter, Marcel A Just, and Peter Shell. What one intelligence test measures: a theoretical account of the processing in the raven progressive matrices test. *Psychological review*, 97(3):404, 1990.
- Raymond Bernard Cattell and Alberta KS Cattell. *Measuring intelligence with the culture fair tests*.
 Institute for Personality and Ability Testing, 1960.
- 577
 578
 578
 578
 579
 579
 580
 579
 580
 579
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
 580
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,
 Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot
 impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*, 2(3):6, 2023.
- Gabriele Chierchia, Delia Fuhrmann, Lisa J Knoll, Blanca Piera Pi-Sunyer, Ashok L Sakhardande,
 and Sarah-Jayne Blakemore. The matrix reasoning item bank (mars-ib): novel, open-access
 abstract reasoning items for adolescents and adults. *Royal Society open science*, 6(10):190232,
 2019.
- François Chollet. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang,
 Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision language models with instruction tuning. *Advances in Neural Information Processing Systems*, 36, 2024.

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

648 649 650	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35: 22199–22213, 2022.
651 652 653	Paulo Guirro Laurence and Elizeu Coutinho Macedo. Cognitive strategies in matrix-reasoning tasks: State of the art. <i>Psychonomic Bulletin & Review</i> , 30(1):147–159, 2023.
654 655 656	Adam Lerer, Sam Gross, and Rob Fergus. Learning physical intuition of block towers by example. In <i>International conference on machine learning</i> , pp. 430–438. PMLR, 2016.
657 658 659	Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge. https://llava-vl.github. io/blog/2024-01-30-llava-next, 2024a.
660 661 662	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. Advances in neural information processing systems, 36, 2024b.
663 664	Nan Liu, Shuang Li, Yilun Du, Josh Tenenbaum, and Antonio Torralba. Learning to compose visual relations. <i>Advances in Neural Information Processing Systems</i> , 34:23166–23178, 2021.
665 666 667 668	Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? <i>arXiv preprint arXiv:2307.06281</i> , 2023.
669 670 671 672	Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. <i>Advances in Neural Information Processing Systems</i> , 35:2507–2521, 2022.
673 674 675	Mikołaj Małkiński and Jacek Mańdziuk. Deep learning methods for abstract visual reasoning: A survey on raven's progressive matrices. <i>arXiv preprint arXiv:2201.12382</i> , 2022.
676 677	Mikołaj Małkiński and Jacek Mańdziuk. A review of emerging research directions in abstract visual reasoning. <i>Information Fusion</i> , 91:713–736, 2023.
678 679 680 681	Mikołaj Małkiński and Jacek Mańdziuk. One self-configurable model to solve many abstract visual reasoning problems. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 14297–14305, 2024.
682 683	AI Meta. Introducing meta llama 3: The most capable openly available llm to date, 2024. URL https://ai. meta. com/blog/meta-llama-3/. Accessed on April, 26, 2024a.
684 685 686	AI Meta. Introducing llama 3.2. URL https://github.com/meta-llama/llama- models/tree/main/models/llama3_2 Accessed on Sep, 2024b.
687 688	Shanka Subhra Mondal, Jonathan D Cohen, and Taylor W Webb. Slot abstractors: Toward scalable abstract visual reasoning. <i>arXiv preprint arXiv:2403.03458</i> , 2024.
689 690 691 692	Arsenii Kirillovich Moskvichev, Victor Vikram Odouard, and Melanie Mitchell. The conceptarc benchmark: Evaluating understanding and generalization in the arc domain. <i>Transactions on Machine Learning Research</i> , 2023.
693 694	Bjorn Ommer and Joachim M Buhmann. Learning the compositional nature of visual objects. In 2007 IEEE Conference on Computer Vision and Pattern Recognition, pp. 1–8. IEEE, 2007.
695 696 697	OpenAI.Gpt-4v(ision) system card.https://openai.com/research/gpt-4v-system-card, 2023.
698 699	OpenAI. Hello gpt-4o. https://openai.com/index/hello-gpt-4o, 2024.
700 701	Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. <i>arXiv preprint arXiv:2306.14824</i> , 2023.

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

756 757 758 750	Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. <i>arXiv</i> preprint arXiv:2311.03079, 2023.
760 761	Taylor Webb, Keith J Holyoak, and Hongjing Lu. Emergent analogical reasoning in large language models. <i>Nature Human Behaviour</i> , 7(9):1526–1541, 2023.
762 763 764	Taylor Webb, Shanka Subhra Mondal, and Jonathan D Cohen. Systematic visual reasoning through object-centric relational abstraction. <i>Advances in Neural Information Processing Systems</i> , 36, 2024a.
765 766 767 768	Taylor W Webb, Steven M Frankland, Awni Altabaa, Simon Segert, Kamesh Krishnamurthy, Declan Campbell, Jacob Russin, Tyler Giallanza, Randall O'Reilly, John Lafferty, et al. The relational bottleneck as an inductive bias for efficient abstraction. <i>Trends in Cognitive Sciences</i> , 2024b.
769 770	Taylor Whittington Webb, Ishan Sinha, and Jonathan Cohen. Emergent symbols through binding in external memory. In <i>International Conference on Learning Representations</i> , 2020.
771 772 773	David Wechsler and Habuku Kodama. <i>Wechsler intelligence scale for children</i> , volume 1. Psycholog- ical corporation New York, 1949.
774 775 776	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. <i>arXiv preprint arXiv:2206.07682</i> , 2022.
777 778 779	Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms. <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , 2023.
780 781 782 783	Jingyi Xu, Tushar Vaidya, Yufei Wu, Saket Chandra, Zhangsheng Lai, and Kai Fong Ernest Chong. Abstract visual reasoning: An algebraic approach for solving raven's progressive matrices. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 6715–6724, 2023a.
784 785 786	Yudong Xu, Wenhao Li, Pashootan Vaezipoor, Scott Sanner, and Elias B Khalil. Llms and the abstrac- tion and reasoning corpus: Successes, failures, and the importance of object-based representations. <i>arXiv preprint arXiv:2305.18354</i> , 2023b.
787 788 789 790	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. <i>arXiv preprint arXiv:2407.10671</i> , 2024.
791 792 793	Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Li- juan Wang. The dawn of lmms: Preliminary explorations with gpt-4v (ision). <i>arXiv preprint</i> <i>arXiv:2309.17421</i> , 9(1):1, 2023.
794 795 796 797	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. <i>Advances in Neural</i> <i>Information Processing Systems</i> , 36, 2024.
798 799 800	Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, et al. Yi: Open foundation models by 01. ai. <i>arXiv preprint arXiv:2403.04652</i> , 2024.
801 802 803	Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it? In <i>The Eleventh International Conference on Learning Representations</i> , 2022.
804 805 806 807	Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. From recognition to cognition: Visual commonsense reasoning. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 6720–6731, 2019.
808 809	Aimen Zerroug, Mohit Vaishnav, Julien Colin, Sebastian Musslick, and Thomas Serre. A benchmark for compositional visual reasoning. <i>Advances in neural information processing systems</i> , 35: 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 nattern recognition, pp. 5317–5327, 2019.
813	

- Baoan Zhang, Junming Yang, Hanjia Lyu, Zijian Jin, Yuan Yao, Mingkai Chen, and Jiebo Luo. Cocot:
 Contrastive chain-of-thought prompting for large multimodal models with multiple image inputs.
 arXiv preprint arXiv:2401.02582, 2024a.
- Buzhen Zhang, Yahan Yu, Chenxing Li, Jiahua Dong, Dan Su, Chenhui Chu, and Dong Yu. Mm-Ilms:
 Recent advances in multimodal large language models. *arXiv preprint arXiv:2401.13601*, 2024b.
- Yizhe Zhang, He Bai, Ruixiang Zhang, Jiatao Gu, Shuangfei Zhai, Josh Susskind, and Navdeep Jaitly. How far are we from intelligent visual deductive reasoning? *arXiv preprint arXiv:2403.04732*, 2024c.
- Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal
 chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023.
- Kai Zhao, Chang Xu, and Bailu Si. Learning visual abstract reasoning through dual-stream networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 16979–16988, 2024.
- Liang Zhou, Kevin A Smith, Joshua B Tenenbaum, and Tobias Gerstenberg. Mental jenga: A counterfactual simulation model of causal judgments about physical support. *Journal of Experimental Psychology: General*, 152(8):2237, 2023.
- Qiji Zhou, Ruochen Zhou, Zike Hu, Panzhong Lu, Siyang Gao, and Yue Zhang. Image-of-thought prompting for visual reasoning refinement in multimodal large language models. *arXiv preprint arXiv:2405.13872*, 2024.
- Beyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.

Appendices

0	
CON	FENTS.

1	Introduction	-
•		,
2	Kelated Works	•
3	MaRs_VAA Datasat	
5	Mans- V QA Dataset	-
4	Visual Cognition Benchmark (VCog-Bench)	4
	4.1 Multi-Image Reasoning via Chain-of-Thought (CoT)	4
	4.2 Vision-Language Models (VLMs)	-
5	Experiments	(
	5.1 Experimental Settings	(
	5.2 Experimental Results	-
	5.2 Ablation Study	ç
		(
		5
	5.5 Visualization	1(
6	Discussion	1(
7	Conclusion	1(
Ap	pendices	17
A	Datasets & Benchmarking Code	18
B	Data Collection and Licenses	18
С	Experimental Settings	19
	C 1 Implementation Details	10
	C 2 System Prompts	20
		20
D	Further Discussion on Limitations and Future Work	20
Е	Ethics Discussion	24
	E.1 Negative Societal Impacts	24
	E.2 Mitigating Bias and Negative Societal Impacts	24
	C	-

918 A DATASETS & BENCHMARKING CODE

920 921 We release the data and annotations of MaRs-VQA anonymously:

- 922 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

925 926 927

928

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.

941

942 CVR The original dataset link of CVR is github.com/serre-lab/CVR. It is under Apache License
943 2.0 (CVR LICENSE). CVR is an accepted paper by NeurIPS 2022 Datasets and Benchmarks track,
944 so all of its data is free to use by public. We follow the same data generation pipeline in CVR to
945 generate 309 samples.

946

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

The paper and study results are under MIT License. All questionnaires are under AttributionNonCommercial 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.

956 After we download all questionnaires from MaRs-IB, we use two Python scripts to merge all question-957 option pairs from different questionnaires into the same sample set. Then, we generate Option Set A, 958 Option Set B in Figure 6 by manipulating the size and image position of option images. After that, 959 we annotate the language description of 4 options in 10 samples from the raw data. The language 960 description is used as system prompt to guide GPT-40 to generate all description for all data in 961 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 962 example of each type of option. 963

⁹⁶⁴ The sub-task statistics of MaRs-VQA is in Table.

965
 966
 967
 968
 969
 969
 969
 969
 960
 961
 962
 963
 964
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965
 965

- 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.

С **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 stimul: 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.

IMPLEMENTATION DETAILS C.1

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.

1032

1034

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052 1053

1054

1055 1056

1058 1059

1061

1062

1063

1064

1067

1068

1069

1070

1071

1074 1075

1033 C.2 System Prompts

For each dataset, we prepare custom system prompt. Their pipeline is similar. First, we created a system message prompt (see Figure 7, 8 for zero-shot inference, and Figure 9, 10, 11 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 $K = \{[r, a, o, s] | r \in \mathcal{R}, a \in \mathcal{A}, o \in \mathcal{O}, s \in S\}$. The output format is a json structure including "Answer" and "Explanation" as keys.

System Prompt for zero-shot inference for MaRs-VQA

System Message

You are a helpful visual reasoning assistant that can solve abstract 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 option images is the correct answer. To select the correct missing shape, you have to deduce relationships between the shapes of the matrix. These shape characteristics varied along these dimensions: shape, color, size, and position in the matrix. You should only respond in the format as described below:

Response Format

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

Figure 7: System prompts for zero-shot MLLM inference of MaRs-VQA.

System Prompt for zero-shot inference for RAVEN

System Message

You are a helpful visual reasoning assistant solve abstract reasoning problem. 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 image is the correct answer. To select the correct missing shape, you have to deduce relationships between the shapes of the matrix. These shape characteristics varied along five dimensions: number, shape (triangle, square, pentagon, hexagon, circle), color (five colors from white to black), size (five size from small to large) and position in the matrix.

You should only respond in the format as described below:

Response Format

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

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

1077 1078

1079 D FURTHER DISCUSSION ON LIMITATIONS AND FUTURE WORK

1080	System Prompt for MLLMs with CoT for MaRs-VOA
1081	
1082	System Message
1083	You are a helpful visual reasoning assistant that can solve abstract visual reasoning problems.
1084	Each task consisted of a question image with a 3 times 3 matrix. Eight of the nine resulting
1085	cells contained an abstract shape, while one cell on the bottom right-hand side of the matrix
1086	mathematic possible alternatives. One of the options is the correct answer
1087	The first step is to describe what is the attribute and relationship between each attribute in
1088	each cell of the 3 times 3 question image. The attributes can be number, position, shape.
1089	size, and color. The cell may contain multiple attributes. The relation might be '3 times 3
1090	sub-blocks', 'rotation', 'insideness'.
1091	The second step is to summarize the relation of three patterns in the first row of the question
1092	image, the relation of three patterns in the second row of the question image, the relation of
1093	two patterns in the third row of the question image.
1094	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
1095	ontions?
1096	Finally, infer what are the potential attributes, objects, relations in the missing cell?
1097	You should only respond in the format as described below:
1090	
1100	Response Format
1101	<i>Explanation</i> : The step-by-step reasoning for the answer.
1102	Answer: The index of the correct answer, as a single letter.
1102	
1104	Eigung 0. System moments for MLLM CoT information of MoDe VOA
1105	Figure 9: System prompts for MILLM Col interence of Mars-VQA.
1106	
1107	Insights Unlike other VOA benchmarks, our work approaches the perspective of human visual
1108	cognition—an underexplored domain Based on our experimental results we offer the following
1109	insights for vision researchers:
1110	
1111	• While scaling laws have some applicability to visual cognition tasks, merely increasing
1112	model size and training data is insufficient to achieve human-level performance.
1113	• To demonstrate that VLMs possess strong visual cognitive abilities, it is crucial to evaluate
1114	them on zero-shot inference tasks like matrix reasoning-tasks characterized by simple
1116	visual content but requiring complex reasoning to find the correct answer.
1117	• Unlike other multi-image visual reasoning benchmarks. VCog Ranch affectively highlights
1110	the performance gap between MLLMs and human cognition in these tasks
1110	the performance sup convert millions and number cognition in these tasks.
1120	From our main and ablation experiments, we observed that as task difficulty increases, the per-
1121	formance of MLLMs in multi-image reasoning scenarios deteriorates. Interestingly, providing
1122	language-based descriptions of each option (i.e., inputting the model with a single question image and
1123	context-based options) improved the models' performance compared to using multi-image options.
1124	This suggests that language still plays a significant role in the visual reasoning processes of current
1125	MLLMs and VLMs.
1126	In contrast, human visual cognition—especially in children—allows individuals to solve matrix
1127	reasoning tasks without relying on advanced language reasoning capabilities. Children can often
1128	solve these tasks effectively by utilizing their visual working memory and pattern recognition skills.
1129	One notential reason for the performance gap is that current MLLMs/VLMs may underemphasize the
1130	visual encoder relative to the language encoder. In many recently released VLMs, the visual module
1131	is much smaller than the language model module, and the visual encoders are frozen during Large
1132	Language Model (LLM) and alignment layer fine-tuning in open-sourced VLMs. This imbalance
1133	might limit the models' capacity to retain and process complex visual information during reasoning
	tasks

1	
	Sustan Massage
	System Message
	Fach task consisted of a question image with a 3 times 3 matrix. Fight 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 ontion 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
	<i>Explanation:</i> The step-by-step reasoning for the answer.
	Answer: The index of the correct answer, as a single letter.
	System Prompt for MLLMs with CoT for CVR
	System Prompt for MLLMs with CoT for CVR
	System Prompt for MLLMs with CoT for CVR System Message You are a helpful visual reasoning assistant that can solve abstract reasoning problems. Each
	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
	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
	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
	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
	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: share: color: size, relationship of
	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
	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
	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.
	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
	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?
	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?
	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:
	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:
	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:
	 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.
	 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: Explanation: The step-by-step reasoning for the answer. Answer: The index of the correct answer, as a single letter.
	 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 <i>Explanation:</i> The step-by-step reasoning for the answer. <i>Answer:</i> The index of the correct answer, as a single letter.
	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 <i>Explanation:</i> The step-by-step reasoning for the answer. <i>Answer:</i> The index of the correct answer, as a single letter. Eigure 11: System prompts for CoT MLLM inference of CVB
	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 <i>Explanation:</i> The step-by-step reasoning for the answer. <i>Answer:</i> The index of the correct answer, as a single letter. Figure 11: System prompts for CoT MLLM inference of CVR.
	 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: Explanation: The step-by-step reasoning for the answer. Answer: The index of the correct answer, as a single letter.
	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.
	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: Figure 11: System prompts for CoT MLLM inference of CVR. better retain visual information during the reasoning process, MLLMs may require more cap and here on bond be complex using process, MLLMs may require more cap and using the reasoning process, MLLMs may require more cap and using the retain visual information during the reasoning process, MLLMs may require more cap and using the redsoning process, MLLMs may require more cap and using the redsoning process, MLLMs may require more cap and using the redsoning process.
	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 o 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: Figure 11: System prompts for CoT MLLM inference of CVR. better retain visual information during the reasoning process, MLLMs may require more can used modules that can handle complex visual patterns and maintain this information througe to reasoning the reasoning process with end to and multiple can be used to and multiple can be and the order or using the relation of the order or using the reasoning to be an using the data and head to and multiple can be applex visual patterns and maintain this information througe to reasoning the the training more using the reasoning the reasoning ton the oreastern with end to and multiple can using the

ine reasoning steps. Moreover, optimizing the training process with end-to-end multimodal training—without freezing any layers in the visual modules—can be beneficial. Recent models have

begun to explore end-to-end VLM fine-tuning, demonstrating the potential of this approach, though
 challenges remain such as the need for multi-round alignment. In the future, developing more
 advanced methods to effectively integrate visual and linguistic features will be crucial.

1191

1192 1193

1194

Limitations In the main paper, we briefly discussed the limitations of our work. Here, we provide 1195 a more in-depth discussion. First, our dataset is composed of limited publicly available matrix 1196 reasoning datasets, which must include human study results. The RAVEN and CVR datasets, created 1197 by the AI/ML community, were not developed following rigorous psychological research norms. 1198 Consequently, our benchmarking results, which utilize these datasets, should not be used to derive 1199 psychological or clinical conclusions. While MaRs-VQA addresses this problem, its samples cannot 1200 represent all formats of matrix reasoning found in IQ tests such as the WISC and the Cattell Culture 1201 Fair Intelligence Test (Cattell & Cattell, 1960). We cannot use these IQ tests directly because they are not freely available, and copyright restrictions usually prevent these pen-and-paper tasks from 1202 being adapted into computerized formats. 1203

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

- 1211
- 1212
- 1213

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

1221 From our experimental results, we observe that current MLLMs have enhanced basic matrix reasoning 1222 capabilities, with models like GPT-40 and Gemini Pro 1.5 achieving significantly higher accuracy than random guessing across all three matrix reasoning tasks. By using Monte Carlo Tree Search 1223 to optimize the results via multi-round reasoning and exclusion, GPT-40 can achieve much better 1224 outcomes, albeit at the cost of increased inference time. We anticipate that the next generation 1225 of MLLMs will approach human-level performance in matrix reasoning. It is crucial to maintain 1226 these visual cognition-based benchmarks, continuously monitor the performance of newly released 1227 MLLMs, and encourage open-source MLLMs and VLMs to include matrix reasoning tasks for 1228 performance comparison. 1229

Finally, we pose the open-ended question of whether MLLMs need to achieve or surpass human-level 1230 zero-shot inference capability in matrix reasoning tasks. Addressing this issue requires drawing 1231 on theories from cognitive science and psychology to understand the nature of human and MLLM 1232 intelligence. Matrix reasoning ability develops early in human neurodevelopment, with children as 1233 young as four providing sensible answers to simple matrix reasoning questions without additional 1234 training, making it a critical component of IQ tests. In contrast, LLMs and MLLMs rely on training 1235 data, fundamentally differing from how children develop cognitive abilities. However, we believe that 1236 these two learning processes share commonalities: both involve the gradual accumulation of skills 1237 and the ability to generalize from past experiences. Exploring these parallels can provide valuable insights into designing MLLMs that more closely mimic human visual cognition, ultimately leading 1239 to more advanced and capable models. Additionally, we observe that current open-source models achieve matrix reasoning performance very close to that of closed-source models. However, VLMs 1240 face challenges in supporting multiple images as input and managing visual memory. Addressing 1241 these challenges is a crucial direction for building more robust open-source VLMs in the future.

1242 E ETHICS DISCUSSION

1243

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

- 1250
- 1251 E.1 NEGATIVE SOCIETAL IMPACTS

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

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

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

1267 1268 E.2 MITIGATING BIAS AND NEGATIVE SOCIETAL IMPACTS

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