

# LOGICVISTA: MULTIMODAL LLM LOGICAL REASONING BENCHMARK IN VISUAL CONTEXTS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We propose LogicVista, an evaluation benchmark that examines multimodal large language models' (MLLMs) integrated **Logical** reasoning capacities in **Visual** contexts. Recent advancements in MLLMs have demonstrated various fascinating abilities such as crafting poetry based on an image to engaging in mathematical reasoning. Despite these feats, there remains a gap in the systematic examination of MLLMs' proficiency in logical reasoning tasks. These skills are routinely invoked in navigation, puzzle-solving, etc. Thus we present LogicVista, which evaluates general logical cognition abilities across a spectrum of 5 logical reasoning tasks with 3 broad capabilities and 11 specific capabilities through a sample of 448 multiple-choice questions. Each is annotated with not only the correct answer but also the human-written reasoning behind the selection, allowing for rich open-ended evaluation as well as MCQ evaluation. A total of 11 MLLMs undergo comprehensive evaluation using LogicVista. We are also introducing a crowdsourced annotation tool to further scale LogicVista with support from the community. Code and Data Available at <https://anonymous.4open.science/r/LogicVista>.

## 1 INTRODUCTION

Recent advancements in Large Language Models (LLMs) are gradually turning the vision of a generalist AI agent into reality. These models exhibit near-human expert-level performance across a variety of tasks and have recently been augmented with visual understanding capabilities, enabling them to tackle even more complex visual challenges. This branch of work, led by both proprietary projects such as Flamingo (Alayrac et al., 2022) and GPT-4 (OpenAI et al., 2024) and open-source works such as LLaVA (Liu et al., 2023a), Mini-GPT4 (Zhu et al., 2023) enhances existing LLMs by incorporating visual comprehension. These models, known as Multimodal Large Language Models (MLLMs), utilize LLMs as the foundation for processing information and generating reasoned outcomes (Yin et al., 2023), bridging the gap between language and vision. Recent MLLMs have demonstrated a range of such intriguing abilities, such as writing poems based on an image (Fu et al., 2023a), engaging in mathematical reasoning (Alayrac et al., 2022), and even aiding in medical diagnosis (Zhang et al., 2023).

**Challenges** Many works have benchmarked MLLM's performance on common multimodal tasks such as recognizing objects (Antol et al., 2015), understanding the text in an image (Singh et al., 2019a), or performing math (Yu et al., 2023; Lu et al., 2024). However, there are two major concerns with existing benchmarks: lack of evaluation of explicit logical-visual-language reasoning skills and potential data leakage in benchmarking data.

Evaluating explicit logical-visual-language reasoning is essential, as it reflects a key aspect of human creativity and intelligence. Proficiency in reasoning skills is widely recognized as a reliable indicator of cognitive ability across various domains (Kahneman, 2012; Träff et al., 2019). While many datasets and benchmarks have been designed to assess the logical reasoning capabilities of AI agents, most are limited to text-based formats, leaving visual reasoning largely underexplored (Liu et al., 2023; Xu et al., 2023; Yang et al., 2023; Lin, 2024; Yang et al., 2024). While some datasets like GLoRE, MathVista, MM-vet, and RAVEN (liu et al., 2023; Lu et al., 2024; Yu et al., 2023; Zhang et al., 2019) have explored aspects of visual logical reasoning, they focus primarily on specific tasks such as mathematical reasoning, spatial reasoning, or world knowledge retrieval, with logical reasoning only partially embedded and not directly analyzed. General-purpose visual question answering and

captioning datasets like TextVQA and VQAv2 (Goyal et al., 2017a; Singh et al., 2019a) contain even fewer examples of visual logical reasoning, concentrating instead on the recognition and identification of visual details. Similarly, specialized benchmarks such as MMMU and OlympiadBench (Yue et al., 2024; He et al., 2024) focus on academic domain questions in subjects like math, science, or history, without directly evaluating the visual logical reasoning capabilities of modern MLLMs.

Moreover, many existing benchmarks rely on publicly available data from the internet, which can easily be included in the training datasets of various models due to its low friction for scraping (as demonstrated in Appendix K). This increases the likelihood that many benchmarking samples are inadvertently leaked into the training data, leading to unfair comparisons of models that do not effectively isolate their reasoning capabilities. In Appendix A, we provide a more comprehensive overview of the gaps in the current literature on MLLM benchmarks.

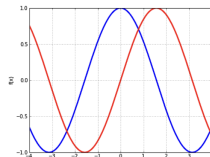
Recently, MLLMs have demonstrated impressive problem-solving and understanding capabilities across various domains. Researchers have aimed to strengthen these models' logical reasoning abilities through novel pre-training techniques, such as directly embedding logical reasoning, as demonstrated with IDOL (Xu et al., 2023). However, their capacity for explicit visual logical reasoning remains largely untested in a comprehensive, systematic way. Thus, developing a scalable and thorough benchmark to assess MLLMs' visual logical reasoning abilities is essential. This would drive advancements in logical reasoning systems within visual contexts, especially as VQA agents gain traction in fields like robotics, biology, and software engineering (Muennighoff et al., 2024; Hong et al., 2023; Xiao et al., 2024), while also providing a framework to evaluate progress in visual understanding and reasoning in MLLMs.



Q: Is the girl touching the ground?

A: No

Reasoning Skill: None  
Capability: Recognition



Q: Which function is monotonic in range  $[0, \pi]$ ?

A: the blue one

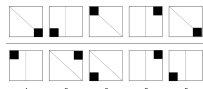
Reasoning Skill: Numerical  
Capability: OCR



Q: What will the girl on the right write?

A: 14

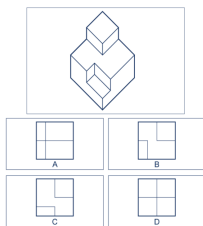
Reasoning Skill: Numerical  
Capability: OCR



Q: Which of the boxes comes next?

A: E

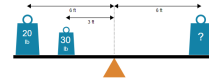
Reasoning Skill: Inductive  
Capability: Diagram



Q: Which of these are the top view?

A: B

Reasoning Skill: Spatial  
Capability: 3D Shape



Q: What is the weight if balanced?

A: C

Reasoning Skill: Mechanical  
Capability: Physics

Figure 1: Capabilities and reasoning skills of different existing benchmarks. The top row shows examples from VQAv2, MathVista, and MM-vet in order from left to right, while the bottom row contains examples from our LogicVista. Unlike previous benchmarks, LogicVista focuses on visual reasoning capacities explicitly.

**This Work** With these motivations, we propose a comprehensive benchmark for general visual logical reasoning to address these challenges. Our benchmark utilizes rigorously sourced data to ensure high quality and fair evaluation of the explicit visual-logical reasoning skills of current state-of-the-art MLLMs. We argue that an effective universal evaluation benchmark should have the following characteristics: (1) coverage of a broad range of human logical reasoning skills, including deductive, inductive, numeric, spatial, and mechanical reasoning; (2) presentation of information in various formats such as Optical Character Recognition (OCR), graphs, charts, and flow diagrams to accommodate diverse data inputs; (3) responses structured for convenient quantitative

analysis, enabling rigorous assessment and comparison of model performance; and (4) scalability to accommodate community feedback and growth, ensuring sustainability and effective evaluation of the benchmark for future models and formats.

To this end, we collect a comprehensive MLLM evaluation benchmark, named LogicVista, which fulfills all the criteria:

- LogicVista covers the examination of 5 representative categories of logical reasoning tasks: inductive (*sample* = 107), deductive (*sample* = 93), numerical (*sample* = 95), spatial (*sample* = 79), and mechanical (*sample* = 74).
- LogicVista covers 3 broad capabilities and 11 specific capabilities to give a comprehensive view of how well MLLMs reason with various visual formats.
- All images- instructions-solution-reasoning are rigorously manually annotated and validated using our robust annotation pipeline.
- Benefiting from our instruction design “please select from A, B, C, D, and E.” and our LLM answer evaluator, we can evaluate different reasoning skills and capabilities and easily perform quantitative statistical analysis based on the natural language output of MLLMs. We also provide more in-depth human-written explanations for why each answer is correct for more through open-ended evaluation.
- To ensure the scalability and sustainability of LogicVista for future evaluations, we introduce the annotation tool used for community crowdsourcing, as detailed in Appendix L.

As shown in Figure 3, LogicVista covers a broad range of reasoning skills, evaluating both open- and closed-source SOTA MLLMs. For example, the question “*Which of these images is the top view of the given object*” in Figure 1(b) requires spatial reasoning from a different perspective, not just object recognition. Since these questions and diagrams are presented without real-world context, they test the MLLM’s core reasoning abilities.

We perform comprehensive evaluations on 11 representative open- and closed-source MLLMs, using 448 samples across 5 key logical reasoning categories, providing the first in-depth assessment of visual logical reasoning in state-of-the-art models like GPT-4 Omni, Claude-3.5 Sonnet, and Gemini-Pro. LogicVista’s evaluation framework breaks down each model’s performance by reasoning skill and capability, offering more nuanced insights than a single overall score. We employ two evaluation methods: MCQ for quick assessments and open-ended chain-of-thought (CoT) for a deeper analysis of the reasoning process, identifying where models succeed or fall short. This approach offers a clearer understanding of each model’s strengths and weaknesses.

Our findings indicate that LogicVista is a highly challenging benchmark, with top-performing models averaging 65% in deductive reasoning but scoring below 30% in other reasoning categories. Notably, GPT-4o and Claude 3.5 Sonnet exhibit state-of-the-art performance on LogicVista, as detailed in Table 2. We observe that most models struggle the most with inductive, numerical, and spatial reasoning, while performing better in deductive and mechanical visual reasoning tasks. Additionally, our analysis shows that MLLMs achieve higher accuracy with MCQ prompts compared to CoT-based prompts, suggesting that MCQs rely more on educated guesses and require less in-depth reasoning. In contrast, CoT prompts often lead to incorrect explanations and lower performance, as models struggle with reasoning or generate hallucinated answers. This pattern reflects human behavior, where selecting a single answer is generally easier than providing a detailed explanation.

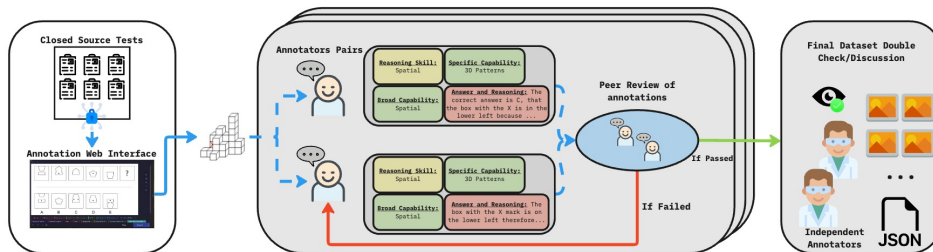
## 2 THE LOGICVISTA DATASET

### 2.1 DATA SOURCES

Most multimodal evaluation benchmarks source images from the open internet, which risks data leakage into MLLM training datasets, potentially giving some models an unfair advantage. To ensure the integrity of LogicVista’s evaluations, we prevent data leakage by collecting and annotating our samples from licensed IQ test sources, with permission from the test creators. These tests are not publicly accessible and require either payment or registration, significantly reducing the likelihood that LogicVista’s samples have been seen by MLLMs during training. Licensing details and sources are also included in the dataset annotations. Additionally, we conducted Google reverse image

162 searches on 50 randomly sampled LogicVista data points, finding that our dataset is not available  
 163 online, whereas many samples from existing benchmarks can be found on the open internet (detailed  
 164 in Appendix K).

## 166 2.2 ANNOTATION AND DATA COLLECTION GUIDELINES



179 Figure 2: LogicVista’s robust manual annotation pipeline ensures high-quality data through multiple  
 180 rounds of peer review and validation.

182 **Motivation** LogicVista comprises images designed to assess the underlying reasoning capabilities  
 183 of MLLMs. Real-life scenes can complicate explicit tests of logical reasoning because they often  
 184 provide contextual clues that enable an AI agent to deduce answers without engaging in direct  
 185 reasoning about the scene. To address this, LogicVista features multiple-choice questions across  
 186 3 broad capabilities and 11 specific capabilities, clearly specifying the type of reasoning required  
 187 without the additional context of real-life scenarios. Such question formats are commonly found  
 188 in intelligence and reasoning tests. Consequently, we initially reviewed over 50 intelligence test  
 189 distributors for suitable tests and formats, focusing on a diverse range of reasoning categories and test  
 190 sizes. This process led us to filter down to approximately 10 closed-source test banks, from which we  
 191 gathered our datasets, seeking permission from the test creators to use their materials for our project.

192 **Annotation Process** To ensure high-quality annotations, we established a rigorous data collection and  
 193 annotation pipeline involving six annotators and two project leads, all of whom are STEM students,  
 194 as detailed in Figure 2. The annotators were organized into pairs, each responsible for annotating the  
 195 same batch of images. They classified each image based on its logical reasoning, broad capability,  
 196 and specific capability, while also providing the correct answer and open-ended reasoning annotations.  
 197 Using an answer key as a reference, annotators developed in-depth explanations for why each answer  
 198 choice was correct.

199 To maintain accuracy and consistency in the open-ended reasoning annotations, the teams collaborated  
 200 to reach a consensus on the correct answers and reasoning for each sample in LogicVista. After  
 201 each annotation sprint, the teams conducted peer reviews, exchanging and refining their annotations.  
 202 Suggested edits were merged into a single batch for each group, which was then submitted to project  
 203 advisors who acted as independent reviewers to ensure the quality of the open-ended reasoning  
 204 annotations and correct answers. Each batch underwent cross-validation by an independent group  
 205 of annotators, providing an objective quality check before incorporation into the final LogicVista  
 206 dataset.

207 At the end of the project, the group reconvened to verify the robustness of all samples, ensuring that  
 208 key annotations, such as open-ended reasoning and question classifications, were both accurate and  
 209 comprehensive. This meticulous process spanned approximately four months. All data were collected  
 210 and annotated from closed sources requiring payment or registration for access, significantly reducing  
 211 the likelihood of the dataset being included in prior training or benchmarking datasets, as outlined in  
 212 Appendix K.

213 **Annotation Categories** To enable a thorough analysis of MLLM performance on visual logical  
 214 reasoning tasks, we provided fine-grained data annotations that allow for examination across various  
 215 aspects. With this goal in mind, we annotated each sample in LogicVista with the following details:  
 the question, the answer, the correct MCQ answer, an open-ended reasoning explanation for why

the MCQ answer is correct, the reasoning skill used, the broad multimodal capability, the specific multimodal capability, and licensing/data source information.

**Annotation Tool** Additionally, we developed an annotation tool, detailed in Appendix L, which we will release for crowdsourcing. This will allow us to scale the pipeline to the broader community, ensuring the sustainability and scalability of LogicVista for future developments.

### 2.3 LOGICVISTA ANALYSIS

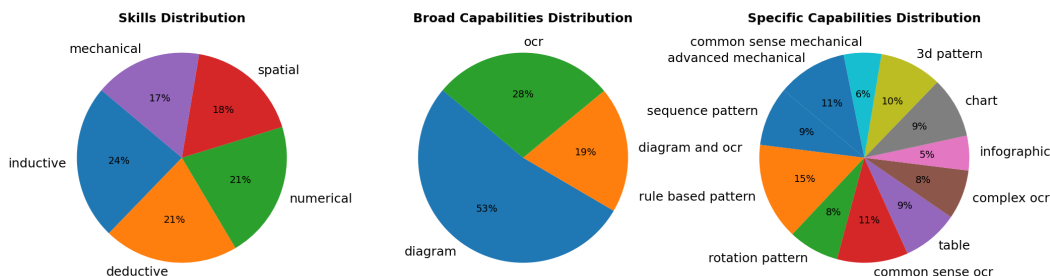


Figure 3: Proportion of reasoning skills and capabilities. On the left is the proportion of questions belonging to each reasoning skill. In the middle is the proportion of questions belonging to broad visual capabilities. On the right is the proportion of questions belonging to specific visual capabilities

A detailed breakdown of the contents in LogicVista is shown in Figure 3. The dataset encompasses 5 core reasoning skills based on fundamental human reasoning capabilities, which we further categorize into broad multimodal capability formats and specific formats for in-depth analysis of MLLM performance in visual logical reasoning. The data is sourced from over 15 human intelligence tests. Samples from the dataset, presented in Appendices I and J, illustrate the richness and diversity of the logical reasoning skills and formats included in LogicVista.

**Multi-modal Capabilities** We define multi-modal capabilities as distinct from reasoning skills, as these capabilities are essential for understanding a multi-modal scene and extracting relevant information. Capabilities refer to the modes in which logical reasoning questions are presented. To ensure comprehensive coverage in LogicVista, we have established a diverse array of 3 *broad capabilities* and 11 *specific capabilities* for evaluation. This division into broad and specific categories provides hierarchical insights into how well MLLMs perform in areas such as OCR versus diagrams at a broader level, while also offering detailed insights into their performance across specific categories, including various diagram presentation styles and formats. This diversity ensures that LogicVista thoroughly evaluates a wide range of logical situations that an MLLM may encounter in everyday reasoning, providing in-depth insights into each capability. Figure 3 illustrates that LogicVista incorporates a balanced mix of various capabilities, including samples that leverage both to solve problems effectively. We define these capabilities in detail in Appendix B.

**Visual Logical Reasoning Skills** The reasoning skills that were of interest for this benchmark were based on common reasoning skills humans use for critical thinking and problem-solving in most contexts derived from popular human intelligence tests. For our evaluation, we summarize these to include the following 5 skills. As seen in Figure 3, LogicVista contains a wide, balanced range of all core reasoning skills. We define these skills in detail in Appendix B.

## 3 EXPERIMENTS

### 3.1 EVALUATION MODELS

To evaluate the performance of MLLMs on LogicVista, we selected a range of representative models detailed in Appendix C Table. 3. Specifically, we selected 11 models for evaluation, including LLaVA (Liu et al., 2023a; 2024), MiniGPT4 (Zhu et al., 2023), Otter (Li et al., 2023a), variations of OpenAI’s GPT-4 (OpenAI et al., 2024), variations of Anthropic’s Claude (Anthropic, 2023),

270 variations of Google’s Gemini (Team et al., 2024), BLIP-2 (Li et al., 2023b), and InstructBLIP (Dai  
271 et al., 2023) We also specifically included models pix2struct (Lee et al., 2023) as they have been  
272 tuned to understand chart or diagram data.

273 We selected a diverse set of models that represent the current MLLM landscape in both open  
274 and closed-source MLLMs. This selection encompasses various model sizes and architectures,  
275 incorporating different visual encoders, backbone language models, and training datasets. Our goal  
276 was to obtain a comprehensive understanding of MLLMs’ performance in visual logical reasoning  
277 skills. The breakdown of models we selected for our experiments is detailed in Table 3.

278 Additionally, we incorporated baseline comparisons to provide a reference for interpreting the results  
279 from MLLMs. First, we established a random baseline that selects answer choices by randomly  
280 sampling from a Gaussian distribution. We also included a frequentist baseline, which selects the  
281 most commonly seen option in the dataset as the response.

### 283 3.2 EVALUATION PROTOCOLS

284 We evaluate the models on LogicVista using two setups: **MCQ-based** prompting assessed with an  
285 LLM-based answer choice extractor, and **CoT-based** prompting evaluated by an LLM-as-judge. We  
286 chose an MCQ-only evaluation strategy for its straightforward and efficient approach to gauging  
287 MLLM performance. The binary nature of MCQ answers (correct/incorrect) simplifies grading and  
288 allows for easy comparisons across various tasks and datasets. This method is also used by several  
289 other datasets, such as MathVista (Lu et al., 2024), establishing its reliability.

290 However, we recognize that MCQ-only evaluations have limitations, as they obscure the reasoning  
291 processes of MLLMs by reducing the evaluation to a binary output without revealing the rationale  
292 behind the answers. To address this, we also incorporate a chain-of-thought evaluation format, where  
293 we ask an LLM judge (GPT-4o) to analyze CoT responses from MLLMs. This judge compares  
294 these responses to the ground truth and explains which aspects were incorrect, providing a finer  
295 understanding of whether MLLMs arrive at the correct answer with sound reasoning or if they  
296 produce incorrect answers despite valid reasoning.

297 To calculate accuracy scores for each model, we use different methods depending on whether we are  
298 evaluating with the MCQ or CoT approach. For MCQ, an LLM-based extractor isolates the selected  
299 answers from the MLLMs’ outputs (which are often full sentences rather than single letters) and  
300 compares them to the correct answers. In the CoT approach, an LLM judge assesses the open-ended  
301 responses against the ground truth. In both cases, the overall logical reasoning score is determined by  
302 dividing the number of correct responses by the total number of samples in that particular category,  
303 whether it pertains to reasoning skills or capabilities.

## 306 4 EXPERIMENTAL RESULTS

### 308 4.1 VISUAL LOGICAL REASONING PERFORMANCE

309 Table 2 highlights the results for these models across five logical reasoning categories. We analyzed  
310 models of different architectures and sizes, comparing them against random and frequentist baselines.

311 Our results show that most models struggle with inductive, numerical, and spatial reasoning, while  
312 generally performing well in deductive and mechanical reasoning tasks.

313 **Training Limitations:** We believe this disparity arises from the limited exposure visual encoders  
314 like CLIP (Dosovitskiy et al., 2021; Radford et al., 2021) have to inductive, numerical, and spatial  
315 reasoning scenarios in their training data. These encoders are typically trained on standard computer  
316 vision (CV) datasets focused on object recognition, classification, and segmentation using text labels.  
317 While this equips models to excel in tasks like identifying and labeling objects or understanding  
318 cause-effect relationships, it leaves them ill-prepared for reasoning on spatial dynamics or inductive  
319 patterns.

320 For instance, LLaVA models, often fine-tuned with data capturing object names and coordinates,  
321 show stronger spatial, inductive, and deductive reasoning than other open-source counterparts. This  
322 underscores the need for vision encoders that capture detailed image information. Despite the  
323

capabilities of advanced backbone LLMs, MLLM reasoning is limited when visual encoders cannot extract crucial logical information. The narrow focus of CV datasets further constrains these models' ability to handle more complex reasoning tasks.

Table 1: LogicVista evaluation results on various multimodal LLMs on broad multi-modal capabilities. Higher scoring models are highlighted **green** and lower scoring models are highlighted **yellow**.

Model	Broad Capabilities			Model	Broad Capabilities		
	Diagram	OCR	Diagram and OCR		Diagram	OCR	Diagram and OCR
Frequentist	26.69%	23.20%	21.84%	pix2struct	9.60%	6.60%	5.75%
Random	22.46%	24.80%	22.99%	miniGPTvicuna7B	11.15%	9.43%	6.90%
Claude 3.5 Sonnet	36.02%	62.40%	39.08%	miniGPTvicuna13B	13.00%	16.98%	12.64%
Claude 3 Opus	30.51%	40.80%	28.74%	instructBLIP-vicuna-7B	12.07%	20.28%	17.24%
Claude 3 Sonnet	30.08%	48.80%	29.89%	instructBLIP-vicuna-13B	10.53%	13.21%	14.94%
Claude 3 Haiku	27.12%	40.80%	16.09%	instructBLIP-flan-t5-xl	20.74%	21.70%	17.24%
GPT4	26.63%	38.68%	25.29%	instructBLIP-flan-t5-xxl	20.12%	25.47%	18.39%
GPT-4o	33.47%	47.20%	26.44%	BLIP2	19.50%	23.11%	18.39%
GPT-4o-mini	25.85%	47.20%	25.29%	LLAVA7B	29.72%	27.36%	26.44%
Gemini-Pro	37.29%	54.40%	32.18%	LLAVA13B	21.67%	24.06%	14.94%
Gemini-Flash	34.75%	45.60%	24.14%	LLAVANEXT-7B-vicuna	26.01%	23.11%	19.54%
otter9B	23.22%	22.17%	18.39%	LLAVANEXT-13B-vicuna	24.15%	23.58%	20.69%

**Architectural Limitations:** Inductive reasoning often involves identifying patterns across multiple examples, which is not emphasized in standard visual training. In contrast, deductive reasoning—grounded in logical structures and patterns common in textual data—is a strength for LLMs due to their extensive training on large text corpora. Numerical reasoning, another area of weakness, requires an understanding of mathematical principles visually—something multi-modal models struggle to integrate with both visual and textual information. Additionally, the architecture of these models may favor certain reasoning types over others. For instance, while attention mechanisms excel at sequential deduction, they may struggle to effectively capture visual spatial relationships. Ultimately, these challenges in reasoning tasks arise from both the limitations in training data and the architectural design of multimodal LLMs. We further elaborate on these points in Section 4.5.

Table 2: LogicVista evaluation results on various multimodal LLMs on each logical reasoning skill. The higher scoring models are highlighted **green** and lower scoring models are highlighted **yellow**.

Model	Logical Reasoning Skills				
	Inductive	Deductive	Numerical	Spatial	Mechanical
Frequentist	25.23%	19.35%	27.37%	26.58%	25.67%
Random	21.50%	30.11%	16.84%	18.99%	29.73%
Claude 3.5 Sonnet	27.10%	65.59%	47.37%	29.11%	52.70%
Claude 3 Opus	21.50%	49.46%	26.32%	25.33%	45.95%
Claude 3 Sonnet	28.04%	53.76%	32.63%	27.85%	33.78%
Claude 3 Haiku	24.30%	47.31%	15.79%	24.05%	33.78%
GPT4	23.36%	54.84%	24.21%	21.52%	41.89%
GPT-4o	23.36%	58.06%	26.32%	26.58%	48.65%
GPT-4o-mini	22.43%	53.76%	26.32%	21.52%	35.14%
Gemini-Pro	28.97%	62.37%	32.63%	24.05%	60.81%
Gemini-Flash	32.71%	51.61%	25.26%	20.25%	50.00%
otter9B	31.78%	24.73%	18.95%	18.99%	21.62%
pix2struct	12.15%	6.45%	2.11%	7.59%	17.57%
miniGPTvicuna7B	10.28%	9.68%	7.37%	3.80%	27.03%
miniGPTvicuna13B	13.08%	23.66%	10.53%	10.13%	17.57%
instructBLIP-vicuna-7B	4.67%	21.51%	24.21%	2.53%	22.97%
instructBLIP-vicuna-13B	3.74%	10.75%	18.95%	5.06%	17.57%
instructBLIP-flan-t5-xl	23.36%	22.58%	22.11%	7.59%	33.78%
instructBLIP-flan-t5-xxl	17.76%	30.11%	24.21%	20.25%	22.97%
BLIP2	17.76%	23.66%	23.16%	24.05%	18.92%
LLAVA7B	29.91%	29.03%	26.32%	25.32%	36.49%
LLAVA13B	18.69%	31.18%	20.00%	27.85%	24.32%
LLAVANEXT-7B-vicuna	26.17%	21.51%	25.26%	27.85%	29.73%
LLAVANEXT-13B-vicuna	22.43%	22.58%	26.32%	26.58%	25.68%
LLAVANEXT-7B-mistral	16.82%	34.41%	23.16%	21.52%	22.97%
LLAVANEXT-34B-NH	20.56%	52.69%	30.53%	24.05%	40.54%

## 4.2 VISUAL CAPABILITIES PERFORMANCE

We highlight the performance of MLLMs on various broad and specific visual capabilities in Appendix J, Tables 1, 4, and 5.

**Broad Capabilities** Our results show that, on average, most models perform better on OCR-type questions than on diagram-format questions within the broad capability category.

A possible reason why multi-modal LLMs (MLLMs) perform better on OCR-type questions compared to diagram-based questions is the difference in visual reasoning complexity. OCR tasks mainly involve recognizing and extracting textual information from images, which plays to the strengths of visual encoders in object recognition and classification. The text in OCR tasks is typically structured, with clear boundaries and minimal need for spatial or abstract reasoning. This allows the MLLM to focus on straightforward text recognition, followed by reasoning using the LLM backbone, which has been shown to excel at various textual reasoning tasks (liu et al., 2023; OpenAI et al., 2024; Touvron et al., 2023). As a result, the multi-modal reasoning task is simplified into a more manageable text-based reasoning process.

In contrast, diagram-based questions typically demand more complex spatial reasoning, pattern recognition, and an understanding of relationships between visual elements. These tasks go beyond merely recognizing objects or labels, requiring the ability to interpret how objects interact, and their relative positions, and sometimes even apply inductive or deductive reasoning. Visual encoders, often not optimized for spatial or abstract relationships, tend to struggle with these challenges. The complexity of interpreting geometric shapes, spatial arrangements, and abstract concepts in diagrams is much greater than the more straightforward task of recognizing and interpreting text in OCR scenarios, as it requires more than simple recognition and identification of basic relationships.

**Specific Capabilities** We found that MLLMs generally perform well on tasks involving complex OCR, common sense OCR, advanced mechanical reasoning, common sense mechanical reasoning, and rotational patterns. However, they tend to struggle with tasks that require understanding 3D patterns and sequential patterns.

This reinforces our earlier hypothesis that MLLMs excel in OCR and mechanical reasoning tasks because these visual formats primarily focus on recognizing simple relationships and identifying objects, rather than interpreting complex spatial interactions. Mechanical formats often depict real-life scenes, making it easier to discern relationships compared to abstract patterns, where the spatial relationships in 3D and sequential formats are more challenging to extract.

In contrast, tasks involving 3D and sequential pattern recognition require a more nuanced understanding of spatial relationships, movement, and order—capabilities that may be underdeveloped in these models due to limitations in their training data and architectures. Spatial and sequential diagram-based tasks, as well as 3D reasoning, demand an advanced ability to comprehend spatial hierarchies and continuous pattern changes—areas where current visual encoders typically struggle. This lack of spatial depth and temporal awareness contributes to the weaknesses observed in MLLMs when addressing more complex reasoning scenarios.

## 4.3 CASE STUDIES ON LOGICVISTA

Our case studies (Figure 4, 9, 10) show that these errors often occur because MLLMs overlook important details or hallucinate facts, yet still guess the correct solution. This underscores the need for better visual encoders that can capture intricate spatial details beyond recognition. In our MiniGPT-4 case study (Figure 4), while the model reaches the correct answer, the left-hand example reveals a failure to grasp key spatial relationships, guessing "C" simply because the question mark is unfilled. This likely stems from the visual encoder’s limitations in interpreting intricate spatial details. Conversely, in the right-hand example, hallucinations lead to incorrect reasoning. Similarly, MiniGPT-4 fabricates details about pipe sizes, resulting in inaccurate reasoning despite correctly identifying certain image elements. Closed-source flagship models also suffer from these visual encoder limitations, as seen in our SOTA model case studies in Appendix F. We also conduct a more in-depth case study analysis of vision encoder performance of MLLMs in Appendix F.

## 4.4 FINE-GRAINED ANSWER ANALYSIS



432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

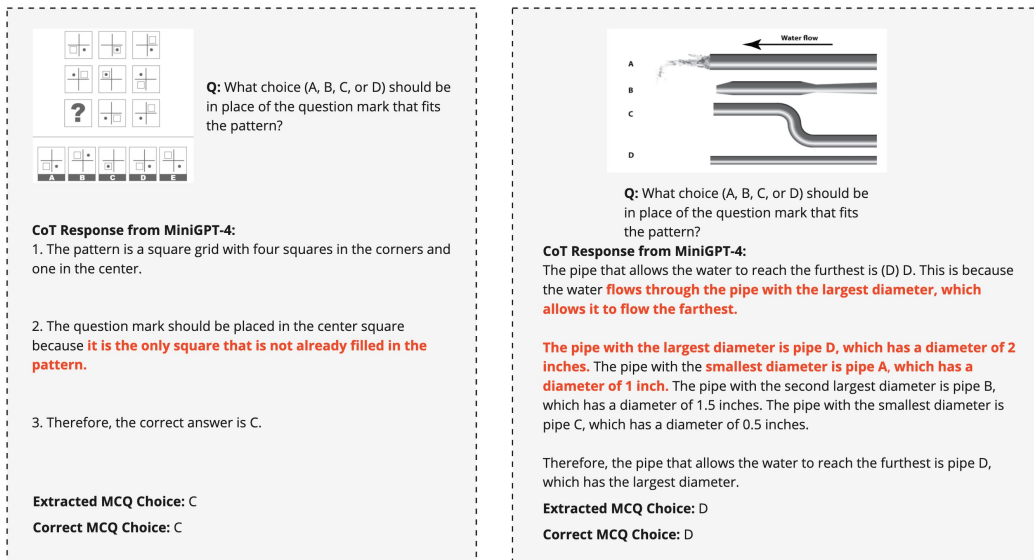


Figure 4: Case study of MiniGPT-4 outputs shows both examples providing correct answers but with flawed CoT reasoning. On the left, the model fails to capture the spatial pattern—where the dot stays constant and the square shifts right—leading to answer C. On the right, the model hallucinates facts about the pipes, missing the key detail that narrower pipes result in faster exit velocity, making D the correct choice.

	Confusion Matrix Analysis (Percentage)			
Correct Answer	37.95	29.69	7.59	15.62
	29.02	22.54	9.60	9.82
Wrong Answer	0.00	0.67	54.46	54.02
	0.00	0.89	61.38	66.74
	Correct Reason		Wrong Reason	

Figure 5: Confusion matrix for performances of SOTA flagship models: Claude 3.5 Sonnet, Gemini Pro, GPT-4o, GPT-4o-mini, arranged clockwise from the top left in each cell.

Sonnet achieves the best results. We observed several instances where models gave correct answers but incorrect explanations. We perform an additional analysis of how MCQ-based evaluation affects evaluation accuracy compared to CoT-based evaluations in Appendix G.

#### 4.5 MODEL COMPARISONS

**Vision Component:** In our evaluation, we considered only open-source vision models for benchmarking. The primary vision encoders used were CLIP-ViT (428M) and EVA-ViT-G (1.13B). When paired

Using chain-of-thought (CoT) evaluations, where an LLM acts as the judge to compare MLLM outputs to ground truth reasoning, we find that most incorrect responses from MLLMs arise from both incorrect answers and flawed reasoning. This suggests the models either fail to fully understand the problem or miss critical details needed to answer accurately. Interestingly, many correct answers still exhibit faulty reasoning, as shown in Figure 5. For example, Claude 3.5 Sonnet, the top-performing model across reasoning skills and capabilities, gave incorrect answers with faulty explanations 54% of the time and correctly answered questions with incorrect explanations 7.9% of the time. Other models, such as Gemini-Pro and GPT-4, displayed similar patterns, with Gemini Pro leading in faulty explanations at 16%. A manual review of Claude 3.5 Sonnet’s responses revealed that over 46% of its outputs included some form of hallucination, either about the image contents (e.g., fabricating facts about patterns or sequences) or general knowledge (e.g., physics, deductive reasoning). Overall, GPT-4o-mini performs the worst, frequently providing incorrect answers and explanations, while Claude 3.5

with large language models (LLMs) such as Vicuna 7B and 13B, the LLaVA variants incorporating CLIP-ViT demonstrated superior performance in spatial, deductive, and inductive reasoning tasks compared to InstructBLIP, which utilized the EVA-ViT-G encoder. Despite these observations, it is challenging to declare a definitive superior model for logical reasoning due to the absence of a comprehensive ablation study, which would provide a more thorough analysis of the strengths and weaknesses of each model configuration. We also present a case study of the failures of modern visual encoders in Appendix F.

**Language Modeling Component:** A range of LLMs, including Vicuna, Flan-T5, and LLaMA, were tested by LogicVista to evaluate their performance. With EVA-ViT-G as the vision encoder, the BLIP-2-12B model combined with Flan-T5-XXL outperformed InstructBLIP, which used Vicuna-13B, particularly in spatial reasoning tasks. This suggests that the Flan-T5 model may possess stronger spatial language processing capabilities compared to Vicuna-13B. This observation highlights the potential impact of different LLMs on the effectiveness of multimodal systems. The performance difference indicates that the choice of LLM can significantly affect the effectiveness of multimodal systems in specific reasoning tasks. Flan-T5’s demonstrated strength in spatial reasoning underscores the importance of choosing LLMs that align with specific reasoning capabilities. However, a more detailed analysis of how LLM logical reasoning performance relates to multimodal logical reasoning could provide better insights into how different LLMs impact the overall performance of visual reasoning systems.

**Training Data:** The comparison of training data performance reveals that MiniGPT, with its datasets including CC3m, SBU, LAION-400M, and a custom set of 3500 images, excels particularly in induction tasks. This suggests that MiniGPT’s training data could be highly effective for tasks requiring the model to generalize from specific inductive examples to broader patterns. On the other hand, InstructBLIP’s training data, which encompasses BLIP2 and 26 transformed datasets, shows stronger performance in a broader range of evaluation categories. This indicates that the suitability of training data may vary depending on the specific types of reasoning or tasks. Some of these datasets may have more samples covering specific reasoning tasks causing different datasets to provide distinct advantages for different reasoning skills.

**Closed/Open-Source Models:** The results suggest that closed-source models like GPT, Gemini, and Claude significantly outperform open-source models in deduction and mechanical reasoning, often with double the accuracy. This advantage likely stems from proprietary optimizations, training techniques, model size, or undisclosed data. Additionally, the continuous updates and fine-tuning specific to these models may contribute to their superior performance. However, in numerical, spatial, and inductive reasoning tasks, both open- and closed-source models show similar effectiveness, with accuracy rates between 22% and 31% across leading closed-source models (GPT, Claude, Gemini) and open-source models (13B LLaVA, Yi models, InstructBLIP). While closed-source models excel in deduction and mechanical reasoning, both model types struggle similarly with spatial and inductive reasoning, suggesting the challenges lie more in the fundamental limitations of current MLLM technologies for visual logical reasoning than in proprietary enhancements. Greater transparency and research could clarify these performance differences and inform future advancements in both open-source and closed-source models, potentially bridging the gap in reasoning capabilities.

## 5 CONCLUSION

In this work, we introduce LogicVista, a comprehensive benchmark designed to evaluate MLLM performance on complex visual logical reasoning tasks, covering inductive, deductive, spatial, numerical, and mechanical visual reasoning skills. We assess 11 state-of-the-art open and closed-source MLLMs, offering insights into the current landscape of these models. Our detailed analysis reveals that MLLMs often struggle with intricate spatial and logical details in images, as their visual encoders are typically trained for broad object recognition. This focus leads to failures in tasks that require a deep understanding of spatial relationships, particularly in inductive, spatial, and numerical reasoning. Our fine-grained CoT case study underscores this limitation, showing that MLLMs tend to generalize rather than capture precise spatial details in both abstract and real-life scenes. We also find that MCQ evaluations often overestimate MLLM performance, as they fail to assess reasoning as effectively as CoT methods. Therefore, we propose future benchmarks emphasize open-ended evaluations that assess the reasoning process, not just final answers.

## REFERENCES

- 540  
541  
542 Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra,  
543 Devi Parikh, Stefan Lee, and Peter Anderson. nocraps: novel object captioning at scale. In *2019*  
544 *IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, October 2019. doi:  
545 10.1109/iccv.2019.00904. URL <http://dx.doi.org/10.1109/ICCV.2019.00904>.
- 546 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel  
547 Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan  
548 Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian  
549 Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo  
550 Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language  
551 model for few-shot learning, 2022.
- 552 Anthropic. Claude (oct 8 version), 2023. URL <https://www.anthropic.com>. Large language  
553 model.
- 554  
555 Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence  
556 Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE International*  
557 *Conference on Computer Vision (ICCV)*, December 2015.
- 558  
559 Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollar, and  
560 C. Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server, 2015.
- 561  
562 Cheng-Han Chiang and Hung yi Lee. Can large language models be an alternative to human  
563 evaluations?, 2023.
- 564  
565 Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang,  
566 Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language  
567 models with instruction tuning, 2023.
- 568  
569 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
570 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,  
571 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale,  
572 2021. URL <https://arxiv.org/abs/2010.11929>.
- 573  
574 Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu  
575 Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation  
576 benchmark for multimodal large language models, 2023a.
- 577  
578 Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. Gptscore: Evaluate as you desire, 2023b.
- 579  
580 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V  
581 in VQA matter: Elevating the role of image understanding in Visual Question Answering. In  
582 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017a.
- 583  
584 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa  
585 matter: Elevating the role of image understanding in visual question answering, 2017b.
- 586  
587 Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu,  
588 Xu Han, Yujie Huang, Yuxiang Zhang, Jie Liu, Lei Qi, Zhiyuan Liu, and Maosong Sun. Olympiad-  
589 bench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal  
590 scientific problems, 2024. URL <https://arxiv.org/abs/2402.14008>.
- 591  
592 Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazhen Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan  
593 Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A  
visual language model for gui agents, 2023. URL <https://arxiv.org/abs/2312.08914>.
- Jinsheng Huang, Liang Chen, Taian Guo, Fu Zeng, Yusheng Zhao, Bohan Wu, Ye Yuan, Haozhe Zhao,  
Zhihui Guo, Yichi Zhang, Jingyang Yuan, Wei Ju, Luchen Liu, Tianyu Liu, Baobao Chang, and  
Ming Zhang. Mmevalpro: Calibrating multimodal benchmarks towards trustworthy and efficient  
evaluation, 2024. URL <https://arxiv.org/abs/2407.00468>.

- 594 Yiqiao Jin, Minje Choi, Gaurav Verma, Jindong Wang, and Srijan Kumar. Mm-soc: Benchmarking  
595 multimodal large language models in social media platforms. In *ACL*, 2024.
- 596  
597 Daniel Kahneman. *Thinking, fast and slow*. Penguin, London, 2012. ISBN 9780141033570  
598 0141033576.
- 599 Kenton Lee, Mandar Joshi, Iulia Turc, Hexiang Hu, Fangyu Liu, Julian Eisenschlos, Urvashi Khan-  
600 delwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. Pix2struct: Screenshot parsing as  
601 pretraining for visual language understanding, 2023.
- 602  
603 Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A  
604 multi-modal model with in-context instruction tuning, 2023a.
- 605  
606 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image  
607 pre-training with frozen image encoders and large language models, 2023b.
- 608  
609 Bill Yuchen Lin. ZeroEval: A Unified Framework for Evaluating Language Models, July 2024. URL  
610 <https://github.com/yuchenlin/ZeroEval>.
- 611  
612 Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. Evaluating the logical  
613 reasoning ability of chatgpt and gpt-4, 2023. URL <https://arxiv.org/abs/2304.03439>.
- 614  
615 Hanmeng liu, Zhiyang Teng, Ruoxi Ning, Jian Liu, Qiji Zhou, and Yue Zhang. Glore: Evaluating  
616 logical reasoning of large language models, 2023. URL <https://arxiv.org/abs/2310.09107>.
- 617  
618 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023a.
- 619  
620 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-  
621 next: Improved reasoning, ocr, and world knowledge, January 2024. URL [https://l1ava-v1.  
622 github.io/blog/2024-01-30-llava-next/](https://l1ava-v1.github.io/blog/2024-01-30-llava-next/).
- 623  
624 Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. G-eval: Nlg  
625 evaluation using gpt-4 with better human alignment, 2023b.
- 626  
627 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi  
628 Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your multi-modal model  
629 an all-around player?, 2023c.
- 630  
631 Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng,  
632 Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning  
633 of foundation models in visual contexts, 2024.
- 634  
635 Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual  
636 question answering benchmark requiring external knowledge, 2019.
- 637  
638 Niklas Muennighoff, Qian Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam  
639 Singh, Xiangru Tang, Leandro von Werra, and Shayne Longpre. Octopack: Instruction tuning code  
640 large language models, 2024. URL <https://arxiv.org/abs/2308.07124>.
- 641  
642 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni  
643 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor  
644 Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian,  
645 Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny  
646 Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks,  
647 Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea  
Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen,  
Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung,  
Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch,  
Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty  
Elet, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte,  
Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel  
Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua  
Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike

- 648 Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon  
649 Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne  
650 Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo  
651 Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar,  
652 Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik  
653 Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich,  
654 Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy  
655 Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie  
656 Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini,  
657 Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne,  
658 Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David  
659 Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie  
660 Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély,  
661 Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo  
662 Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano,  
663 Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng,  
664 Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto,  
665 Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power,  
666 Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis  
667 Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted  
668 Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel  
669 Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon  
670 Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,  
671 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang,  
672 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston  
673 Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,  
674 Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason  
675 Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff,  
676 Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu,  
677 Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba,  
678 Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang,  
679 William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024.
- 678 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
679 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.  
680 Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- 682 Leonard Salewski, A. Sophia Koepke, Hendrik P. A. Lensch, and Zeynep Akata. *CLEVR-X: A  
683 Visual Reasoning Dataset for Natural Language Explanations*, pp. 69–88. Springer International  
684 Publishing, 2022. ISBN 9783031040832. doi: 10.1007/978-3-031-04083-2\_5. URL [http://dx.doi.org/10.1007/978-3-031-04083-2\\_5](http://dx.doi.org/10.1007/978-3-031-04083-2_5).
- 686 Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for  
687 image captioning with reading comprehension, 2020.
- 689 Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh,  
690 and Marcus Rohrbach. Towards vqa models that can read, 2019a.
- 692 Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh,  
693 and Marcus Rohrbach. Towards vqa models that can read, 2019b.
- 694 Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut,  
695 Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson,  
696 Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy  
697 Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom  
698 Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli  
699 Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Jack  
700 Krawczyk, Cosmo Du, Ed Chi, Heng-Tze Cheng, Eric Ni, Purvi Shah, Patrick Kane, Betty Chan,  
701 Manaal Faruqui, Aliaksei Severyn, Hanzhao Lin, YaGuang Li, Yong Cheng, Abe Ittycheriah,  
Mahdis Mahdieh, Mia Chen, Pei Sun, Dustin Tran, Sumit Bagri, Balaji Lakshminarayanan,

702 Jeremiah Liu, Andras Orban, Fabian Güra, Hao Zhou, Xinying Song, Aurelien Boffy, Harish  
703 Ganapathy, Steven Zheng, HyunJeong Choe, Ágoston Weisz, Tao Zhu, Yifeng Lu, Siddharth  
704 Gopal, Jarrod Kahn, Maciej Kula, Jeff Pitman, Rushin Shah, Emanuel Taropa, Majd Al Mery,
705 Martin Baeuml, Zhifeng Chen, Laurent El Shafey, Yujing Zhang, Olcan Sercinoglu, George Tucker,  
706 Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs,  
707 Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas  
708 Gonzalez, Misha Khalman, Jakub Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp,  
709 Lev Proleev, Yi Luan, Xi Chen, James Lottes, Nathan Schucher, Federico Lebron, Alban Rustemi,  
710 Natalie Clay, Phil Crone, Tomas Kocisky, Jeffrey Zhao, Bartek Perz, Dian Yu, Heidi Howard, Adam  
711 Bloniarz, Jack W. Rae, Han Lu, Laurent Sifre, Marcello Maggioni, Fred Alcober, Dan Garrette,  
712 Megan Barnes, Shantanu Thakoor, Jacob Austin, Gabriel Barth-Maron, William Wong, Rishabh  
713 Joshi, Rahma Chaabouni, Deeni Fatiha, Arun Ahuja, Gaurav Singh Tomar, Evan Senter, Martin  
714 Chadwick, Ilya Kornakov, Nithya Attaluri, Iñaki Iturrate, Ruibo Liu, Yunxuan Li, Sarah Cogan,  
715 Jeremy Chen, Chao Jia, Chenjie Gu, Qiao Zhang, Jordan Grimstad, Ale Jakse Hartman, Xavier  
716 Garcia, Thanumalayan Sankaranarayana Pillai, Jacob Devlin, Michael Laskin, Diego de Las Casas,  
717 Dasha Valter, Connie Tao, Lorenzo Blanco, Adrià Puigdomènech Badia, David Reitter, Mianna  
718 Chen, Jenny Brennan, Clara Rivera, Sergey Brin, Shariq Iqbal, Gabriela Surita, Jane Labanowski,  
719 Abhi Rao, Stephanie Winkler, Emilio Parisotto, Yiming Gu, Kate Olszewska, Ravi Addanki,  
720 Antoine Miech, Annie Louis, Denis Teplyashin, Geoff Brown, Elliot Catt, Jan Balaguer, Jackie  
721 Xiang, Pidong Wang, Zoe Ashwood, Anton Briukhov, Albert Webson, Sanjay Ganapathy, Smit  
722 Sanghavi, Ajay Kannan, Ming-Wei Chang, Axel Stjerngren, Josip Djolonga, Yuting Sun, Ankur  
723 Bapna, Matthew Aitchison, Pedram Pejman, Henryk Michalewski, Tianhe Yu, Cindy Wang, Juliette  
724 Love, Junwhan Ahn, Dawn Bloxwich, Kehang Han, Peter Humphreys, Thibault Sellam, James  
725 Bradbury, Varun Godbole, Sina Samangooei, Bogdan Damoc, Alex Kaskasoli, Sébastien M. R.  
726 Arnold, Vijay Vasudevan, Shubham Agrawal, Jason Riesa, Dmitry Lepikhin, Richard Tanburn,  
727 Srivatsan Srinivasan, Hyeontaek Lim, Sarah Hodgkinson, Pranav Shyam, Johan Ferret, Steven Hand,  
728 Ankush Garg, Tom Le Paine, Jian Li, Yujia Li, Minh Giang, Alexander Neitz, Zaheer Abbas, Sarah  
729 York, Machel Reid, Elizabeth Cole, Aakanksha Chowdhery, Dipanjan Das, Dominika Rogozińska,  
730 Vitaliy Nikolaev, Pablo Sprechmann, Zachary Nado, Lukas Zilka, Flavien Prost, Luheng He,  
731 Marianne Monteiro, Gaurav Mishra, Chris Welty, Josh Newlan, Dawei Jia, Miltiadis Allamanis,  
732 Clara Huiyi Hu, Raoul de Liedekerke, Justin Gilmer, Carl Saroufim, Shruti Rijhwani, Shaobo Hou,  
733 Disha Shrivastava, Anirudh Baddepudi, Alex Goldin, Adnan Ozturel, Albin Cassirer, Yunhan Xu,  
734 Daniel Sohn, Devendra Sachan, Reinald Kim Amplayo, Craig Swanson, Dessie Petrova, Shashi  
735 Narayan, Arthur Guez, Siddhartha Brahma, Jessica Landon, Miteyan Patel, Ruizhe Zhao, Kevin  
736 Villela, Luyu Wang, Wenhao Jia, Matthew Rahtz, Mai Giménez, Legg Yeung, James Keeling,  
737 Petko Georgiev, Diana Mincu, Boxi Wu, Salem Haykal, Rachel Saputro, Kiran Vodrahalli, James  
738 Qin, Zeynep Cankara, Abhanshu Sharma, Nick Fernando, Will Hawkins, Behnam Neyshabur,  
739 Solomon Kim, Adrian Hutter, Priyanka Agrawal, Alex Castro-Ros, George van den Driessche,  
740 Tao Wang, Fan Yang, Shuo yiin Chang, Paul Komarek, Ross McIlroy, Mario Lučić, Guodong  
741 Zhang, Wael Farhan, Michael Sharman, Paul Natsev, Paul Michel, Yamini Bansal, Siyuan Qiao,  
742 Kris Cao, Siamak Shakeri, Christina Butterfield, Justin Chung, Paul Kishan Rubenstein, Shivani  
743 Agrawal, Arthur Mensch, Kedar Soparkar, Karel Lenc, Timothy Chung, Aedon Pope, Loren  
744 Maggione, Jackie Kay, Priya Jhakra, Shibo Wang, Joshua Maynez, Mary Phuong, Taylor Tobin,  
745 Andrea Tacchetti, Maja Trebacz, Kevin Robinson, Yash Katariya, Sebastian Riedel, Paige Bailey,  
746 Kefan Xiao, Nimesh Ghelani, Lora Aroyo, Ambrose Slone, Neil Houlsby, Xuehan Xiong, Zhen  
747 Yang, Elena Gribovskaya, Jonas Adler, Mateo Wirth, Lisa Lee, Music Li, Thais Kagohara, Jay  
748 Pavagadhi, Sophie Bridgers, Anna Bortsova, Sanjay Ghemawat, Zafarali Ahmed, Tianqi Liu,  
749 Richard Powell, Vijay Bolina, Mariko Iinuma, Polina Zablotskaia, James Besley, Da-Woon Chung,  
750 Timothy Dozat, Ramona Comanescu, Xiance Si, Jeremy Greer, Guolong Su, Martin Polacek,  
751 Raphaël Lopez Kaufman, Simon Tokumine, Hexiang Hu, Elena Buchatskaya, Yingjie Miao,  
752 Mohamed Elhawaty, Aditya Siddhant, Nenad Tomasev, Jinwei Xing, Christina Greer, Helen Miller,  
753 Shereen Ashraf, Aurko Roy, Zizhao Zhang, Ada Ma, Angelos Filos, Milos Besta, Rory Blevins,  
754 Ted Klimenko, Chih-Kuan Yeh, Soravit Changpinyo, Jiaqi Mu, Oscar Chang, Mantas Pajarskas,  
755 Carrie Muir, Vered Cohen, Charline Le Lan, Krishna Haridasan, Amit Marathe, Steven Hansen,  
Sholto Douglas, Rajkumar Samuel, Mingqiu Wang, Sophia Austin, Chang Lan, Jiepu Jiang, Justin  
Chiu, Jaime Alonso Lorenzo, Lars Lowe Sjösund, Sébastien Cevey, Zach Gleicher, Thi Avrahami,  
Anudhyan Boral, Hansa Srinivasan, Vittorio Selo, Rhys May, Konstantinos Aisopos, Léonard  
Hussenot, Livio Baldini Soares, Kate Baumli, Michael B. Chang, Adrià Recasens, Ben Caine,  
Alexander Pritzel, Filip Pavetic, Fabio Pardo, Anita Gergely, Justin Frye, Vinay Ramasesh, Dan

756 Horgan, Kartikeya Badola, Nora Kassner, Subhrajit Roy, Ethan Dyer, Víctor Campos Campos, Alex  
757 Tomala, Yunhao Tang, Dalia El Badawy, Elspeth White, Basil Mustafa, Oran Lang, Abhishek Jindal,  
758 Sharad Vikram, Zhitao Gong, Sergi Caelles, Ross Hemsley, Gregory Thornton, Fangxiaoyu Feng,  
759 Wojciech Stokowiec, Ce Zheng, Phoebe Thacker, Çağlar Ünlü, Zhishuai Zhang, Mohammad Saleh,  
760 James Svensson, Max Bileschi, Piyush Patil, Ankesh Anand, Roman Ring, Katerina Tsihlias, Arpi  
761 Vezer, Marco Selvi, Toby Shevlane, Mikel Rodriguez, Tom Kwiatkowski, Samira Daruki, Keran  
762 Rong, Allan Dafoe, Nicholas FitzGerald, Keren Gu-Lemberg, Mina Khan, Lisa Anne Hendricks,  
763 Marie Pellat, Vladimir Feinberg, James Cobon-Kerr, Tara Sainath, Maribeth Rauh, Sayed Hadi  
764 Hashemi, Richard Ives, Yana Hasson, Eric Noland, Yuan Cao, Nathan Byrd, Le Hou, Qingze  
765 Wang, Thibault Sottiaux, Michela Paganini, Jean-Baptiste Lespiau, Alexandre Moufarek, Samer  
766 Hassan, Kaushik Shivakumar, Joost van Amersfoort, Amol Mandhane, Pratik Joshi, Anirudh Goyal,  
767 Matthew Tung, Andrew Brock, Hannah Sheahan, Vedant Misra, Cheng Li, Nemanja Rakićević,  
768 Mostafa Dehghani, Fangyu Liu, Sid Mittal, Junhyuk Oh, Seb Noury, Eren Sezener, Fantine Huot,  
769 Matthew Lamm, Nicola De Cao, Charlie Chen, Sidharth Mudgal, Romina Stella, Kevin Brooks,  
770 Gautam Vasudevan, Chenxi Liu, Mainak Chain, Nivedita Melinkeri, Aaron Cohen, Venus Wang,  
771 Kristie Seymore, Sergey Zubkov, Rahul Goel, Summer Yue, Sai Krishnakumaran, Brian Albert,  
772 Nate Hurley, Motoki Sano, Anhad Mohananey, Jonah Joughin, Egor Filonov, Tomasz Kępa, Yomna  
773 Eldawy, Jiawern Lim, Rahul Rishi, Shirin Badiezadegan, Taylor Bos, Jerry Chang, Sanil Jain, Sri  
774 Gayatri Sundara Padmanabhan, Subha Puttagunta, Kalpesh Krishna, Leslie Baker, Norbert Kalb,  
775 Vamsi Bedapudi, Adam Kurzrok, Shuntong Lei, Anthony Yu, Oren Litvin, Xiang Zhou, Zhichun  
776 Wu, Sam Sobell, Andrea Siciliano, Alan Papir, Robby Neale, Jonas Bragagnolo, Tej Toor, Tina  
777 Chen, Valentin Anklin, Feiran Wang, Richie Feng, Milad Gholami, Kevin Ling, Lijuan Liu, Jules  
778 Walter, Hamid Moghaddam, Arun Kishore, Jakob Adamek, Tyler Mercado, Jonathan Mallinson,  
779 Siddhinita Wandekar, Stephen Cagle, Eran Ofek, Guillermo Garrido, Clemens Lombriser, Maksim  
780 Mukha, Botu Sun, Hafeezul Rahman Mohammad, Josip Matak, Yadi Qian, Vikas Peswani, Pawel  
781 Janus, Quan Yuan, Leif Schelin, Oana David, Ankur Garg, Yifan He, Oleksii Duzhyi, Anton  
782 Älgmyr, Timothée Lottaz, Qi Li, Vikas Yadav, Luyao Xu, Alex Chinien, Rakesh Shivanna,  
783 Aleksandr Chuklin, Josie Li, Carrie Spadine, Travis Wolfe, Kareem Mohamed, Subhabrata Das,  
784 Zihang Dai, Kyle He, Daniel von Dincklage, Shyam Upadhyay, Akanksha Maurya, Luyan Chi,  
785 Sebastian Krause, Khalid Salama, Pam G Rabinovitch, Pavan Kumar Reddy M, Aarush Selvan,  
786 Mikhail Dekhtyarev, Golnaz Ghiasi, Erdem Guven, Himanshu Gupta, Boyi Liu, Deepak Sharma,  
787 Idan Heimlich Shtacher, Shachi Paul, Oscar Akerlund, François-Xavier Aubet, Terry Huang, Chen  
788 Zhu, Eric Zhu, Elico Teixeira, Matthew Fritze, Francesco Bertolini, Liana-Eleonora Marinescu,  
789 Martin Bølle, Dominik Paulus, Khyatti Gupta, Tejasi Latkar, Max Chang, Jason Sanders, Roopa  
790 Wilson, Xuewei Wu, Yi-Xuan Tan, Lam Nguyen Thiet, Tulsee Doshi, Sid Lall, Swaroop Mishra,  
791 Wanming Chen, Thang Luong, Seth Benjamin, Jasmine Lee, Ewa Andrejczuk, Dominik Rabiej,  
792 Vipul Ranjan, Krzysztof Styrz, Pengcheng Yin, Jon Simon, Malcolm Rose Harriott, Mudit Bansal,  
793 Alexei Robsky, Geoff Bacon, David Greene, Daniil Mirylenka, Chen Zhou, Obaid Sarvana,  
794 Abhimanyu Goyal, Samuel Andermatt, Patrick Siegler, Ben Horn, Assaf Israel, Francesco Pongetti,  
795 Chih-Wei "Louis" Chen, Marco Selvatici, Pedro Silva, Kathie Wang, Jackson Tolins, Kelvin Guu,  
796 Roey Yogev, Xiaochen Cai, Alessandro Agostini, Maulik Shah, Hung Nguyen, Noah Ó Donnai,le,  
797 Sébastien Pereira, Linda Friso, Adam Stambler, Adam Kurzrok, Chenkai Kuang, Yan Romanikhin,  
798 Mark Geller, ZJ Yan, Kane Jang, Cheng-Chun Lee, Wojciech Fica, Eric Malmi, Qijun Tan, Dan  
799 Banica, Daniel Balle, Ryan Pham, Yanping Huang, Diana Avram, Hongzhi Shi, Jasjot Singh, Chris  
800 Hidey, Niharika Ahuja, Pranab Saxena, Dan Dooley, Srividya Pranavi Potharaju, Eileen O'Neill,  
801 Anand Gokulchandran, Ryan Foley, Kai Zhao, Mike Dusenberry, Yuan Liu, Pulkit Mehta, Ragha  
802 Kotikalapudi, Chalence Safranek-Shrader, Andrew Goodman, Joshua Kessinger, Eran Globen,  
803 Prateek Kolhar, Chris Gorgolewski, Ali Ibrahim, Yang Song, Ali Eichenbaum, Thomas Brovelli,  
804 Sahitya Potluri, Preethi Lahoti, Cip Baetu, Ali Ghorbani, Charles Chen, Andy Crawford, Shalini  
805 Pal, Mukund Sridhar, Petru Gurita, Asier Mujika, Igor Petrovski, Pierre-Louis Cedoz, Chenmei Li,  
806 Shiyuan Chen, Niccolò Dal Santo, Siddharth Goyal, Jitesh Punjabi, Karthik Kappaganthu, Chester  
807 Kwak, Pallavi LV, Sarmishta Velury, Himadri Choudhury, Jamie Hall, Premal Shah, Ricardo  
808 Figueira, Matt Thomas, Minjie Lu, Ting Zhou, Chintu Kumar, Thomas Jurdi, Sharat Chikkerur,  
809 Yenai Ma, Adams Yu, Soo Kwak, Victor Åhdel, Sujevan Rajayogam, Travis Choma, Fei Liu,  
Aditya Barua, Colin Ji, Ji Ho Park, Vincent Hellendoorn, Alex Bailey, Taylan Bilal, Huanjie Zhou,  
Mehrdad Khatir, Charles Sutton, Wojciech Rzadkowski, Fiona Macintosh, Konstantin Shagin, Paul  
Medina, Chen Liang, Jinjing Zhou, Pararth Shah, Yingying Bi, Attila Dankovics, Shipra Banga,  
Sabine Lehmann, Marissa Bredesen, Zifan Lin, John Eric Hoffmann, Jonathan Lai, Raynald Chung,  
Kai Yang, Nihal Balani, Arthur Bražinskas, Andrei Sozanschi, Matthew Hayes, Héctor Fernández

810 Alcalde, Peter Makarov, Will Chen, Antonio Stella, Liselotte Snijders, Michael Mandl, Ante  
811 Kärroman, Paweł Nowak, Xinyi Wu, Alex Dyck, Krishnan Vaidyanathan, Raghavender R, Jessica  
812 Mallet, Mitch Rudominer, Eric Johnston, Sushil Mittal, Akhil Udathu, Janara Christensen, Vishal  
813 Verma, Zach Irving, Andreas Santucci, Gamaleldin Elsayed, Elnaz Davoodi, Marin Georgiev, Ian  
814 Tenney, Nan Hua, Geoffrey Cideron, Edouard Leurent, Mahmoud Alnahlawi, Ionut Georgescu,  
815 Nan Wei, Ivy Zheng, Dylan Scandinaro, Heinrich Jiang, Jasper Snoek, Mukund Sundararajan,  
816 Xuezi Wang, Zack Ontiveros, Itay Karo, Jeremy Cole, Vinu Rajashekhar, Lara Tumeh, Eyal Ben-  
817 David, Rishub Jain, Jonathan Uesato, Romina Datta, Oskar Bunyan, Shimu Wu, John Zhang, Piotr  
818 Stanczyk, Ye Zhang, David Steiner, Subhajt Naskar, Michael Azzam, Matthew Johnson, Adam  
819 Paszke, Chung-Cheng Chiu, Jaume Sanchez Elias, Afroz Mohiuddin, Faizan Muhammad, Jin  
820 Miao, Andrew Lee, Nino Vieillard, Jane Park, Jiageng Zhang, Jeff Stanway, Drew Garmon, Abhijit  
821 Karmarkar, Zhe Dong, Jong Lee, Aviral Kumar, Luowei Zhou, Jonathan Evens, William Isaac,  
822 Geoffrey Irving, Edward Loper, Michael Fink, Isha Arkatkar, Nanxin Chen, Izhak Shafran, Ivan  
823 Petrychenko, Zhe Chen, Johnson Jia, Anselm Levskaya, Zhenkai Zhu, Peter Grabowski, Yu Mao,  
824 Alberto Magni, Kaisheng Yao, Javier Snaider, Norman Casagrande, Evan Palmer, Paul Suganthan,  
825 Alfonso Castaño, Irene Giannoumis, Wooyeol Kim, Mikołaj Rybiński, Ashwin Sreevatsa, Jennifer  
826 Prendki, David Soergel, Adrian Goedeckemeyer, Willi Gierke, Mohsen Jafari, Meenu Gaba, Jeremy  
827 Wiesner, Diana Gage Wright, Yawen Wei, Harsha Vashisht, Yana Kulizhskaya, Jay Hoover, Maigo  
828 Le, Lu Li, Chimezie Iwuanyanwu, Lu Liu, Kevin Ramirez, Andrey Khorlin, Albert Cui, Tian  
829 LIN, Marcus Wu, Ricardo Aguilar, Keith Pallo, Abhishek Chakladar, Ginger Perng, Elena Allica  
830 Abellan, Mingyang Zhang, Ishita Dasgupta, Nate Kushman, Ivo Penchev, Alena Repina, Xihui Wu,  
831 Tom van der Weide, Priya Ponnappalli, Caroline Kaplan, Jiri Simsa, Shuangfeng Li, Olivier Dousse,  
832 Fan Yang, Jeff Piper, Nathan Ie, Rama Pasumarthi, Nathan Lintz, Anitha Vijayakumar, Daniel  
833 Andor, Pedro Valenzuela, Minnie Lui, Cosmin Paduraru, Daiyi Peng, Katherine Lee, Shuyuan  
834 Zhang, Somer Greene, Duc Dung Nguyen, Paula Kurylowicz, Cassidy Hardin, Lucas Dixon, Lili  
835 Janzer, Kiam Choo, Ziqiang Feng, Biao Zhang, Achintya Singhal, Dayou Du, Dan McKinnon,  
836 Natasha Antropova, Tolga Bolukbasi, Orgad Keller, David Reid, Daniel Finchelstein, Maria Abi  
837 Raad, Remi Crocker, Peter Hawkins, Robert Dadashi, Colin Gaffney, Ken Franko, Anna Bulanova,  
838 Rémi Leblond, Shirley Chung, Harry Askham, Luis C. Cobo, Kelvin Xu, Felix Fischer, Jun Xu,  
839 Christina Sorokin, Chris Alberti, Chu-Cheng Lin, Colin Evans, Alek Dimitriev, Hannah Forbes,  
840 Dylan Banarse, Zora Tung, Mark Omernick, Colton Bishop, Rachel Sterneck, Rohan Jain, Jiawei  
841 Xia, Ehsan Amid, Francesco Piccinno, Xingyu Wang, Praseem Banzal, Daniel J. Mankowitz, Alex  
842 Polozov, Victoria Krakovna, Sasha Brown, MohammadHossein Bateni, Dennis Duan, Vlad Firoiu,  
843 Meghana Thotakuri, Tom Natan, Matthieu Geist, Ser tan Girgin, Hui Li, Jiayu Ye, Ofir Roval,  
844 Reiko Tojo, Michael Kwong, James Lee-Thorp, Christopher Yew, Danila Sinopalnikov, Sabela  
845 Ramos, John Mellor, Abhishek Sharma, Kathy Wu, David Miller, Nicolas Sonnerat, Denis Vnukov,  
846 Rory Greig, Jennifer Beattie, Emily Caveness, Libin Bai, Julian Eisenschlos, Alex Korchemniy,  
847 Tomy Tsai, Mimi Jasarevic, Weize Kong, Phuong Dao, Zeyu Zheng, Frederick Liu, Fan Yang,  
848 Rui Zhu, Tian Huey Teh, Jason Sanmiya, Evgeny Gladchenko, Nejc Trdin, Daniel Toyama, Evan  
849 Rosen, Sasan Tavakkol, Linting Xue, Chen Elkind, Oliver Woodman, John Carpenter, George  
850 Papamakarios, Rupert Kemp, Sushant Kafle, Tanya Grunina, Rishika Sinha, Alice Talbert, Diane  
851 Wu, Denese Owusu-Afriyie, Cosmo Du, Chloe Thornton, Jordi Pont-Tuset, Pradyumna Narayana,  
852 Jing Li, Saaber Fatehi, John Wieting, Omar Ajmeri, Benigno Uria, Yeongil Ko, Laura Knight,  
853 Amélie Héliou, Ning Niu, Shane Gu, Chenxi Pang, Yeqing Li, Nir Levine, Ariel Stolovich, Rebeca  
854 Santamaria-Fernandez, Sonam Goenka, Wenny Yustalim, Robin Strudel, Ali Elqursh, Charlie  
855 Deck, Hyo Lee, Zonglin Li, Kyle Levin, Raphael Hoffmann, Dan Holtmann-Rice, Olivier Bachem,  
856 Sho Arora, Christy Koh, Soheil Hassas Yeganeh, Siim Pöder, Mukarram Tariq, Yanhua Sun,  
857 Lucian Ionita, Mojtaba Seyedhosseini, Pouya Tafti, Zhiyu Liu, Anmol Gulati, Jasmine Liu, Xinyu  
858 Ye, Bart Chrzaszcz, Lily Wang, Nikhil Sethi, Tianrun Li, Ben Brown, Shreya Singh, Wei Fan,  
859 Aaron Parisi, Joe Stanton, Vinod Koverkathu, Christopher A. Choquette-Choo, Yunjie Li, TJ Lu,  
860 Abe Ittycheriah, Prakash Shroff, Mani Varadarajan, Sanaz Bahargam, Rob Willoughby, David  
861 Gaddy, Guillaume Desjardins, Marco Cornero, Brona Robenek, Bhavishya Mittal, Ben Albrecht,  
862 Ashish Shenoy, Fedor Moiseev, Henrik Jacobsson, Alireza Ghaffarkhah, Morgane Rivière, Alanna  
863 Walton, Clément Crepy, Alicia Parrish, Zongwei Zhou, Clement Farabet, Carey Radebaugh,  
Praveen Srinivasan, Claudia van der Salm, Andreas Fildjeland, Salvatore Scellato, Eri Latorre-  
Chimoto, Hanna Klimczak-Plucińska, David Bridson, Dario de Cesare, Tom Hudson, Piermaria  
Mendolicchio, Lexi Walker, Alex Morris, Matthew Mauger, Alexey Guseynov, Alison Reid, Seth  
Odoom, Lucia Loher, Victor Cotruta, Madhavi Yenugula, Dominik Grewe, Anastasia Petrushkina,  
Tom Duerig, Antonio Sanchez, Steve Yadlowsky, Amy Shen, Amir Globerson, Lynette Webb,



- 864 Sahil Dua, Dong Li, Surya Bhupatiraju, Dan Hurt, Haroon Qureshi, Ananth Agarwal, Tomer  
865 Shani, Matan Eyal, Anuj Khare, Shreyas Rammohan Belle, Lei Wang, Chetan Tekur, Mihir Sanjay  
866 Kale, Jinliang Wei, Ruoxin Sang, Brennan Saeta, Tyler Liechty, Yi Sun, Yao Zhao, Stephan  
867 Lee, Pandu Nayak, Doug Fritz, Manish Reddy Vuyyuru, John Aslanides, Nidhi Vyas, Martin  
868 Wicke, Xiao Ma, Evgenii Eltyshev, Nina Martin, Hardie Cate, James Manyika, Keyvan Amiri,  
869 Yelin Kim, Xi Xiong, Kai Kang, Florian Luisier, Nilesh Tripuraneni, David Madras, Mandy Guo,  
870 Austin Waters, Oliver Wang, Joshua Ainslie, Jason Baldridge, Han Zhang, Garima Pruthi, Jakob  
871 Bauer, Feng Yang, Riham Mansour, Jason Gelman, Yang Xu, George Polovets, Ji Liu, Honglong  
872 Cai, Warren Chen, XiangHai Sheng, Emily Xue, Sherjil Ozair, Christof Angermueller, Xiaowei  
873 Li, Anoop Sinha, Weiren Wang, Julia Wiesinger, Emmanouil Koukoumidis, Yuan Tian, Anand  
874 Iyer, Madhu Gurumurthy, Mark Goldenson, Parashar Shah, MK Blake, Hongkun Yu, Anthony  
875 Urbanowicz, Jennimaria Palomaki, Chrisantha Fernando, Ken Durden, Harsh Mehta, Nikola  
876 Momchev, Elahe Rahimtoroghi, Maria Georgaki, Amit Raul, Sebastian Ruder, Morgan Redshaw,  
877 Jinhyuk Lee, Denny Zhou, Komal Jalan, Dinghua Li, Blake Hechtman, Parker Schuh, Milad Nasr,  
878 Kieran Milan, Vladimir Mikulik, Juliana Franco, Tim Green, Nam Nguyen, Joe Kelley, Aroma  
879 Mahendru, Andrea Hu, Joshua Howland, Ben Vargas, Jeffrey Hui, Kshitij Bansal, Vikram Rao,  
880 Rakesh Ghiya, Emma Wang, Ke Ye, Jean Michel Sarr, Melanie Moranski Preston, Madeleine  
881 Elish, Steve Li, Aakash Kaku, Jigar Gupta, Ice Pasupat, Da-Cheng Juan, Milan Someswar, Tejvi  
882 M., Xinyun Chen, Aida Amini, Alex Fabrikant, Eric Chu, Xuanyi Dong, Amruta Muthal, Senaka  
883 Buthpitiya, Sarthak Jauhari, Nan Hua, Urvashi Khandelwal, Ayal Hitron, Jie Ren, Larissa Rinaldi,  
884 Shahar Drath, Avigail Dabush, Nan-Jiang Jiang, Harshal Godhia, Uli Sachs, Anthony Chen,  
885 Yicheng Fan, Hagai Taitelbaum, Hila Noga, Zhuyun Dai, James Wang, Chen Liang, Jenny Hamer,  
886 Chun-Sung Ferng, Chenel Elkind, Aviel Atias, Paulina Lee, Vít Listík, Mathias Carlen, Jan van de  
887 Kerkhof, Marcin Pikus, Krunoslav Zaher, Paul Müller, Sasha Zykova, Richard Stefanec, Vitaly  
888 Gatsko, Christoph Hirsenschall, Ashwin Sethi, Xingyu Federico Xu, Chetan Ahuja, Beth Tsai,  
889 Anca Stefanoiu, Bo Feng, Keshav Dhandhania, Manish Katyal, Akshay Gupta, Atharva Parulekar,  
890 Divya Pitta, Jing Zhao, Vivaan Bhatia, Yashodha Bhavnani, Omar Alhadlaq, Xiaolin Li, Peter  
891 Danenberg, Dennis Tu, Alex Pine, Vera Filippova, Abhipso Ghosh, Ben Limonchik, Bhargava  
892 Urala, Chaitanya Krishna Lanka, Derik Clive, Yi Sun, Edward Li, Hao Wu, Kevin Hongtongsak,  
893 Ianna Li, Kalind Thakkar, Kuanysh Omarov, Kushal Majmudar, Michael Alverson, Michael  
894 Kucharski, Mohak Patel, Mudit Jain, Maksim Zabelin, Paolo Pelagatti, Rohan Kohli, Saurabh  
895 Kumar, Joseph Kim, Swetha Sankar, Vineet Shah, Lakshmi Ramachandruni, Xiangkai Zeng, Ben  
896 Bariach, Laura Weidinger, Tu Vu, Alek Andreev, Antoine He, Kevin Hui, Sheleem Kashem, Amar  
897 Subramanya, Sissie Hsiao, Demis Hassabis, Koray Kavukcuoglu, Adam Sadovsky, Quoc Le,  
898 Trevor Strohman, Yonghui Wu, Slav Petrov, Jeffrey Dean, and Oriol Vinyals. Gemini: A family of  
899 highly capable multimodal models, 2024. URL <https://arxiv.org/abs/2312.11805>.
- 900 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
901 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand  
902 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language  
903 models, 2023.
- 904 Ulf Träff, Linda Olsson, Kenny Skagerlund, Mikael Skagenholt, and Rickard Östergren. Logical  
905 reasoning, spatial processing, and verbal working memory: Longitudinal predictors of physics  
906 achievement at age 12–13 years. *Frontiers in Psychology*, 10, August 2019. ISSN 1664-1078. doi:  
907 10.3389/fpsyg.2019.01929. URL <http://dx.doi.org/10.3389/fpsyg.2019.01929>.
- 908 Yijia Xiao, Edward Sun, Yiqiao Jin, Qifan Wang, and Wei Wang. Proteingpt: Multimodal llm for  
909 protein property prediction and structure understanding, 2024. URL <https://arxiv.org/abs/2408.11363>.
- 910 Zihang Xu, Ziqing Yang, Yiming Cui, and Shijin Wang. Idol: Indicator-oriented logic pre-training  
911 for logical reasoning, 2023. URL <https://arxiv.org/abs/2306.15273>.
- 912 Yuan Yang, Siheng Xiong, Ali Payani, Ehsan Shareghi, and Faramarz Fekri. Harnessing the power  
913 of large language models for natural language to first-order logic translation, 2023. URL <https://arxiv.org/abs/2305.15541>.
- 914 Yuan Yang, Siheng Xiong, Ali Payani, Ehsan Shareghi, and Faramarz Fekri. Can llms reason in the  
915 wild with programs?, 2024.

918 Zhengyuan Yang, Yijuan Lu, Jianfeng Wang, Xi Yin, Dinei Florencio, Lijuan Wang, Cha Zhang, Lei  
919 Zhang, and Jiebo Luo. Tap: Text-aware pre-training for text-vqa and text-caption, 2020.  
920

921 Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on  
922 multimodal large language models, 2023.

923 Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and  
924 Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities, 2023.  
925

926 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu  
927 Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin,  
928 Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhua Chen.  
929 Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert  
930 agi, 2024. URL <https://arxiv.org/abs/2311.16502>.

931 Chi Zhang, Feng Gao, Baoxiong Jia, Yixin Zhu, and Song-Chun Zhu. Raven: A dataset for relational  
932 and analogical visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and  
933 Pattern Recognition (CVPR)*, 2019.

934 Xiaoman Zhang, Chaoyi Wu, Ziheng Zhao, Weixiong Lin, Ya Zhang, Yanfeng Wang, and Weidi Xie.  
935 Pmc-vqa: Visual instruction tuning for medical visual question answering, 2023.  
936

937 Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing  
938 vision-language understanding with advanced large language models, 2023.  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

## Appendix: LogicVista: Multimodal LLM Logical Reasoning Benchmark in Visual Contexts

### LIMITATIONS

A limitation of our work is the absence of a human baseline for comparison. Having this baseline would offer valuable insights into how MLLMs perform relative to humans. Although several of our test sources suggest that average human performance is around 75%, this figure varies across different reasoning skills, which is why we chose not to include it in our study for the sake of reliability and accuracy. A more dependable approach would be to conduct multiple human trials to establish a consistent average performance for comparison.

Additionally, while our dataset size is comparable to other multimodal benchmarks like MM-vet (Yu et al., 2023), it is relatively smaller than some larger-scale benchmarks such as MMBench or MMMU (Liu et al., 2023c; Yue et al., 2024). To address this, we will release a crowdsourcing annotation tool, detailed in Appendix L, to further scale LogicVista in the future.

To address both concerns and promote further research, we have also open-sourced these reasoning annotations. They are now publicly available for the community, providing a valuable resource for training and improving the logical reasoning capabilities of multimodal LLMs. We encourage future work to make full use of these annotations to develop more comprehensive and contextually rich evaluation methods.

### A RELATED WORKS

**LLM-Based Evaluation.** LogicVista adopts an open-ended LLM-based evaluation approach, which facilitates the generation and assessment of diverse answer styles and question types beyond the limitations of binary or multiple-choice responses. This innovative method leverages the capabilities of large language models (LLMs) for comprehensive model evaluation, a technique that has been effectively applied in natural language processing (NLP) tasks and other VQA benchmarks (Chiang & yi Lee, 2023; Liu et al., 2023b; Fu et al., 2023b; Jin et al., 2024; Lu et al., 2024). Our findings show that this LLM-based evaluation framework is both versatile and robust, providing a unified and flexible assessment across different modalities, including open- and closed-ended responses. By accommodating a broad range of answer styles and question types, this approach deepens and expands model evaluation, leading to a more comprehensive understanding of model performance.

**Vision-Language Benchmarks** Traditional vision-language benchmarks have largely focused on evaluating specific perceptual abilities. Datasets like MM-vet, RAVEN, CLEVR-X, and TextVQA each address distinct aspects of visual recognition: TextVQA emphasizes recognition-based VQA, testing how well models can caption and accurately describe key image details; MM-vet evaluates world knowledge, basic math, detail capture, and OCR in recognition tasks and everyday scene reasoning. Meanwhile, RAVEN and CLEVR-X assess spatial relation recognition in 2D and 3D objects, providing insights into how well MLLMs understand spatial reasoning (Goyal et al., 2017b; Yu et al., 2023; Zhang et al., 2019; Singh et al., 2019a; Sidorov et al., 2020; Salewski et al., 2022). Image captioning and description generation have also been extensively studied (Chen et al., 2015; Agrawal et al., 2019), along with more specialized tasks like scene text understanding (Singh et al., 2019b; Sidorov et al., 2020; Yang et al., 2020) and integrating external knowledge (Marino et al., 2019). Other benchmarks, such as OlympiadBench (He et al., 2024), focus on Olympiad-level math and science challenges to compare MLLMs with human performance. Large-scale multidisciplinary benchmarks like MMMU (Yue et al., 2024) assess MLLMs across a range of subjects, including science, math, humanities, and history.

### B DATASET DEFINITIONS

Here we define concretely what each of our capabilities and logical reasoning skill categories refer to.

## 1026 B.1 LOGICAL REASONING SKILLS

1027

1028 We define our 5 **logical reasoning skills** based on common human visual reasoning abilities. Our  
 1029 goal is to assess how effectively MLLMs perform in general reasoning skills that humans rely on for  
 1030 everyday problem-solving. These skills reflect the types of reasoning that MLLM agents are likely to  
 1031 encounter in real-world settings where they may be deployed. Our definitions are largely inspired by  
 1032 traditional human IQ and intelligence tests.

1033

1034 • **Inductive Reasoning** the ability to infer the next entry in a pattern given a pattern of  
 1035 observations. It is the ability to make generalizations based on some observations and  
 1036 make an educated guess. It moves from many specific observations to a generalization. An  
 1037 example could be given observations that when John eats dairy products, he gets a stomach  
 1038 ache. An inductive conclusion can be drawn that he is most likely lactose intolerant.

1038

1039 • **Deductive Reasoning** the ability to conclude a specific case when given a general principle  
 1040 or pattern. It moves from the general to the specific. An example could be given the  
 1041 statement “all men are mortal”, one can conclude that “John is mortal” because John is a  
 1042 man.

1042

1043 • **Numerical Reasoning** the ability to read arithmetic problems in the image and solve the  
 1044 math equations. An example could be given the equation “ $10 + 10 = ?$ ”, the answer would  
 1045 be “20”.

1045

1046 • **Spatial Reasoning** the ability to understand the spatial relationship between objects and  
 1047 patterns and reason with those relationships. An example could be seeing an unfolded box  
 1048 and understanding what the box could look like when it is folded up.

1048

1049 • **Mechanical Reasoning** the ability to recognize a physical system and solve equations based  
 1050 on that system or answer questions about that system. An example could be seeing a set  
 1051 of 3 gears and understanding which gears will turn clockwise and which ones will turn  
 1052 counterclockwise.

1052

## 1053 B.2 BROAD AND SPECIFIC CAPABILITIES

1054

1055 We categorize our multi-modal capabilities into broad and specific classifications to gain hierarchical  
 1056 insights into which information formats are better or worse understood by MLLMs.

1057

1058 Here we present our definitions for **broad capabilities**:

1058

1059 • **Optical Character Recognition (OCR)**: refers to the ability to reason over text inside  
 1060 images and scenes.

1061

1062 • **Diagrams**: refers to the ability to reason about diagrams that represent real-life scenes,  
 1063 abstract logic, spatial relationships, and more.

1063

1064 • **Mixed (Both OCR and Diagram)**: refers to an integration of both OCR and diagrams,  
 1065 where comprehending the text and the visual elements within the image is essential for  
 1066 accurately answering the question.

1066

1067 Here we present our definitions for **specific capabilities**:

1068

1069 • **Chart**: refers to numerical charts and graphs.

1069

1070 • **Infographic**: refers to infographic-style puzzles that illustrate both real-life and abstract  
 1071 scenes.

1071

1072 • **Table**: refers to words and numbers only tables depicting some trend or concept.

1072

1073 • **Common Sense OCR**: refers to text questions describing common everyday situations  
 1074 using common English words.

1074

1075 • **Complex OCR**: refers to text questions describing technical or highly abstract situations  
 1076 using jargon and complex sentences.

1076

1077 • **Rotation Pattern**: Patterns and puzzles that necessitate an understanding of 2D and/or 3D  
 1078 object rotations.

1077

1079 • **3D Pattern**: Patterns that require 3D spatial relation understandings.

- **Rule Based Pattern:** Patterns that require understanding of a set of externally defined rules.
- **Sequence Pattern:** Patterns presented in a strictly sequential format, typically involved with induction.
- **Common Sense Mechanical:** Puzzles concerned with a common sense understanding of basic physics and mechanics.
- **Advanced Mechanical:** Puzzles concerned with an advanced and specialized understanding of physics and mechanics.

## C SELECTED MLLMs FOR EVALUATION

Model	Size	Language Model	Vision Model
Claude 3.5 Sonnet	N/A <sup>1</sup>	N/A	N/A
Claude 3 Opus	N/A	N/A	N/A
Claude 3 Sonnet	N/A	N/A	N/A
Claude 3 Haiku	N/A	N/A	N/A
GPT-4 Vision	N/A	N/A	N/A
GPT-4o	N/A	N/A	N/A
GPT-4o-mini	N/A	N/A	N/A
Gemini Pro	N/A	N/A	N/A
Gemini Flash	N/A	N/A	N/A
Otter-9B	9B	MPT-7B	CLIP ViT-L/14
Pix2Struct	1.3B	ViT	ViT
MiniGPT-4-7B	7B	Vicuna-7B	BLIP-2 Q-Former
MiniGPT-4-13B	13B	Vicuna-13B	BLIP-2 Q-Former
InstructBLIP-Vicuna-7B	7B	Vicuna-7B	BLIP-2 Q-Former
InstructBLIP-Vicuna-13B	13B	Vicuna-13B	BLIP-2 Q-Former
InstructBLIP-FLAN-T5-xl	3B	FLAN-T5 XL	BLIP-2 Q-Former
InstructBLIP-FLAN-T5-xxl	11B	FLAN-T5 XXL	BLIP-2 Q-Former
BLIP-2	2.7B	OPT-2.7B	EVA-ViT-G
LLaVA-Vicuna-7B	7B	Vicuna-7B	CLIP ViT-L/14
LLaVA-Vicuna-13B	13B	Vicuna-13B	CLIP ViT-L/336px
LLaVA-NeXT-Mistral-7B	7B	Mistral-7B	CLIP ViT-L/14
LLaVA-NeXT-Vicuna-7B	7B	Vicuna-7B	CLIP ViT-L/14
LLaVA-NeXT-Vicuna-13B	13B	Vicuna-13B	CLIP ViT-L/336px
LLaVA-NeXT-Nous-Hermes-Yi-34B	34B	Nous Hermes 2-Yi-34B	CLIP ViT-L/336px

Table 3: Summary of the MLLMs used for evaluations in this study. Model details for close-sourced models like Claude, GPT, and Gemini are not open to the public.

## D BROAD AND SPECIFIC VISUAL CAPABILITIES EVALUATION

We present tabular results evaluating various SOTA open-source and closed-source MLLM models in Table 1, 4, and 5, analyzing their performance across different visual capabilities.

## E SOTA MODEL EVALUATION RESULT

We present graphs illustrating the evaluations of key SOTA closed-source flagship models. Our analysis shows that Claude 3.5 Sonnet consistently performs well across all categories of reasoning and capabilities, with GPT-4o and Gemini Pro following closely in second place.

<sup>1</sup>N/A: Not disclosed

Table 4: Model evaluation results on various multimodal LLMs for Specific Capabilities (Part 1). The highest scoring models are highlighted **green** and lower scoring models are highlighted **yellow**.

Model	3D Pattern	Rule Based Pattern	Sequence Pattern	Rotation Pattern	Table	Chart
Frequentist	27.91%	26.87%	17.07%	31.43%	30.77%	23.81%
Random	23.26%	20.90%	21.95%	14.29%	7.69%	23.81%
Claude 3.5 Sonnet	25.58%	32.84%	17.07%	34.29%	51.28%	50.00%
Claude 3 Opus	16.28%	25.37%	21.95%	28.57%	17.95%	33.33%
Claude 3 Sonnet	27.91%	31.34%	21.95%	28.57%	35.90%	40.48%
Claude 3 Haiku	23.26%	22.39%	24.39%	28.57%	25.64%	14.29%
GPT4	25.58%	26.87%	12.20%	20.00%	28.21%	26.19%
GPT-4o	27.91%	28.36%	19.51%	20.00%	17.95%	30.95%
GPT-4o-mini	27.91%	23.88%	19.51%	14.29%	23.08%	33.33%
Gemini-Pro	23.26%	26.87%	31.71%	25.71%	33.33%	38.10%
Gemini-Flash	13.95%	38.81%	21.95%	28.57%	25.64%	33.33%
otter9B	11.63%	37.31%	24.39%	25.71%	28.21%	11.90%
pix2struct	4.65%	7.46%	17.07%	14.29%	7.69%	0.00%
miniGPTvicuna7B	4.65%	8.96%	12.20%	2.86%	10.26%	7.14%
miniGPTvicuna13B	11.63%	14.93%	12.20%	5.71%	10.26%	14.29%
instructBLIP-vicuna-7B	4.65%	5.97%	2.44%	0.00%	23.08%	23.81%
instructBLIP-vicuna-13B	4.65%	4.48%	7.32%	0.00%	15.38%	21.43%
instructBLIP-flan-t5-xl	9.30%	26.87%	4.88%	20.00%	28.21%	19.05%
instructBLIP-flan-t5-xxl	23.26%	19.40%	14.63%	17.14%	33.33%	21.43%
BLIP2	20.93%	17.91%	14.63%	31.43%	28.21%	16.67%
LLAVA7B	27.91%	32.84%	24.39%	22.86%	23.08%	21.43%
LLAVA13B	27.91%	16.42%	24.39%	25.71%	25.64%	19.05%
LLAVANEXT-7B-vicuna	34.88%	26.87%	19.51%	25.71%	30.77%	19.05%
LLAVANEXT-13B-vicuna	27.91%	22.39%	17.07%	31.43%	30.77%	21.43%
LLAVANEXT-7B-mistral	13.95%	14.93%	21.95%	28.57%	20.51%	23.81%
LLAVANEXT-34B-NH	27.91%	19.40%	24.39%	17.14%	28.21%	19.05%

Table 5: Model evaluation results on various multimodal LLMs for Specific Capabilities (Part 2). The highest scoring models are highlighted **green** and lower scoring models are highlighted **yellow**.

Model	Infographic	Complex OCR	Common Sense OCR	Advanced Mechanical	Common Sense Mechanical
Frequentist	20.83%	17.65%	22.45%	25.00%	26.92%
Random	16.67%	35.29%	30.61%	25.00%	38.46%
Claude 3.5 Sonnet	37.50%	73.53%	63.27%	50%	57.69%
Claude 3 Opus	29.17%	58.82%	46.94%	37.50%	64.54%
Claude 3 Sonnet	16.67%	55.88%	55.10%	25.00%	50.00%
Claude 3 Haiku	12.50%	47.06%	48.98%	25.00%	50.00%
GPT4	67.65%	53.06%	31.25%	61.54%	31.25%
GPT-4o	73.08%	35.42%	73.08%	73.08%	73.08%
GPT-4o-mini	52.94%	63.27%	29.17%	46.15%	46.15%
Gemini-Pro	58.82%	69.39%	58.33%	65.38%	65.38%
Gemini-Flash	8.33%	50.00%	59.18%	41.67%	65.38%
otter9B	29.41%	20.41%	22.92%	19.23%	19.23%
pix2struct	8.16%	16.67%	19.23%	19.23%	19.23%
miniGPTvicuna7B	31.25%	19.23%	19.23%	19.23%	19.23%
miniGPTvicuna13B	32.35%	20.41%	16.67%	19.23%	19.23%
instructBLIP-vicuna-7B	32.35%	16.33%	25.00%	19.23%	19.23%
instructBLIP-vicuna-13B	17.65%	6.12%	25.00%	19.23%	19.23%
instructBLIP-flan-t5-xl	32.35%	16.33%	25.00%	19.23%	50.00%
instructBLIP-flan-t5-xxl	29.41%	28.57%	25.00%	19.23%	19.23%
BLIP2	17.65%	30.61%	22.92%	19.23%	11.54%
LLAVA7B	31.25%	26.53%	31.25%	46.15%	46.15%
LLAVA13B	18.75%	34.62%	18.75%	34.62%	34.62%
LLAVANEXT-7B-vicuna	23.53%	22.45%	27.08%	34.62%	34.62%
LLAVANEXT-13B-vicuna	20.83%	20.41%	20.83%	34.62%	34.62%
LLAVANEXT-7B-mistral	31.25%	18.75%	30.77%	34.62%	30.77%
LLAVANEXT-34B-NH	55.88%	59.18%	41.67%	38.46%	38.46%

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241



Figure 6: SOTA evaluation results of CoT evaluations on logical reasoning skills. As seen here, Claude 3.5 Sonnet has superior performance.

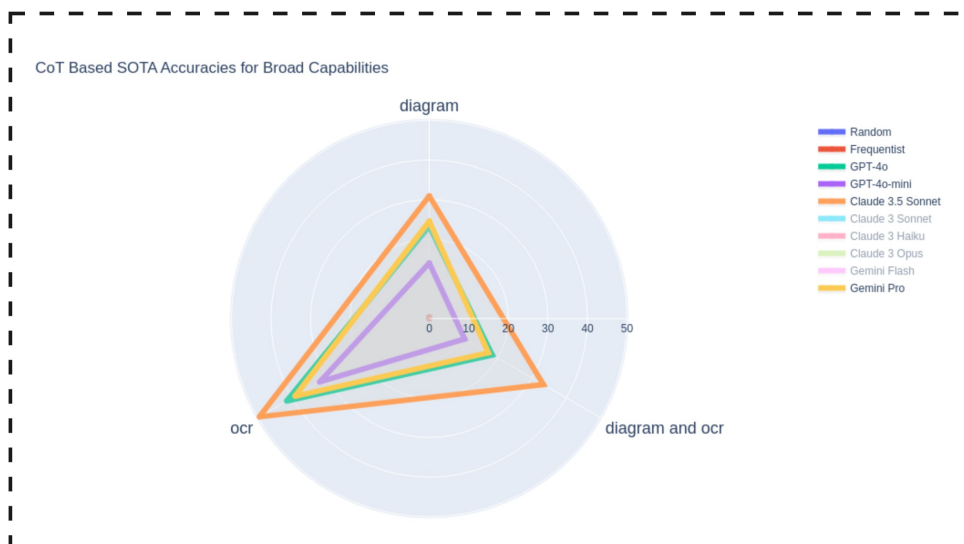


Figure 7: SOTA evaluation results of CoT evaluations on broad visual capabilities. As seen here, Claude 3.5 Sonnet has superior performance.

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

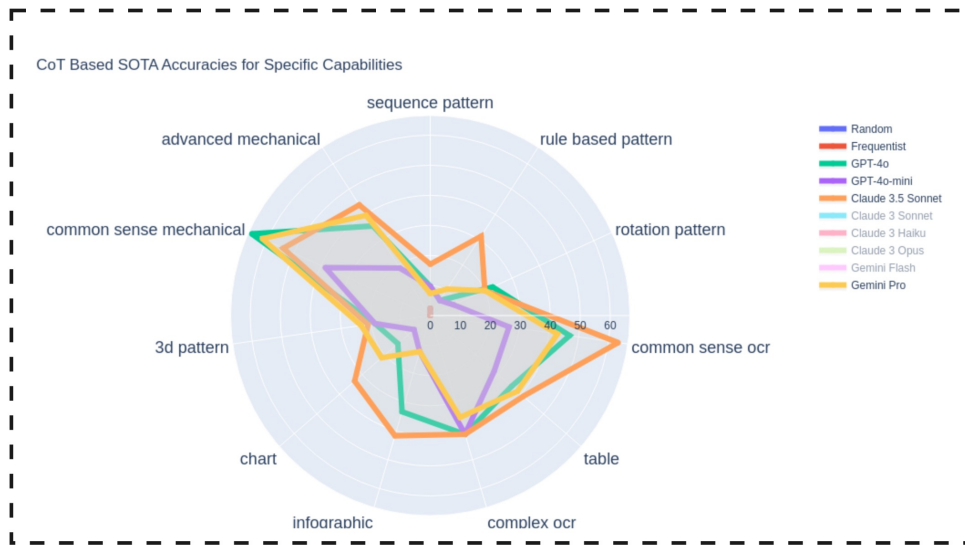


Figure 8: SOTA evaluation results of CoT evaluations on specific visual capabilities. As seen here, Claude 3.5 Sonnet has superior performance. However, it is bested by GPT-4o in some categories like common sense mechanical formats and complex OCR.

## F SOTA CLOSED-SOURCE MLLMS CASE STUDIES

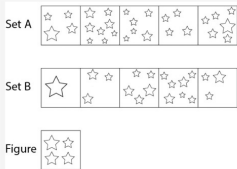
Given that Claude 3.5 Sonnet is the top-performing MLLM across multiple of LogicVista’s question types and formats, we conducted a case study to examine its shortcomings in capturing the spatial and logical relationships essential for complex visual reasoning.

In Figures 9 and 10, Claude 3.5 Sonnet struggles to identify key spatial relationships, such as the shape and number of stars, while also overlooking the logical context necessary for answering LogicVista questions in a CoT format. Similarly, in Figure 10, Claude misinterprets the sequential movement between the circle and triangle. These examples highlight a common issue with modern MLLM vision encoders: they tend to focus on object recognition rather than understanding the relationships between objects, which is essential for accurate visual logical reasoning.

Building on this insight, we conducted an additional experiment to assess the specific details a modern SOTA MLLM can capture. Using Claude 3.5 Sonnet, the top-performing model for LogicVista, we prompted it to provide detailed descriptions of LogicVista samples with a focus on spatial relationships between objects. As shown in Figure 11, MLLMs excel at simpler recognition tasks, such as identifying spring lengths, enabling the model to solve the problem easily. However, when tasked with recognizing more complex spatial relationships, current MLLMs struggle. For instance, in Figure 12, Claude misses intricate spatial patterns and instead focuses on broad features—reflecting a limitation of traditional CV encoders, which are good at general visual recognition but struggle to accurately capture specific spatial arrangements like the positioning of hexagons, squares, and circles. Even in less abstract cases, such as the one on the right in Figure 12 depicting a sled, Claude fails to distinguish key details like the width and size of runners, instead hallucinating differences in sled sizes. When asked specifically about the runner sizes, Claude either misidentifies them as similar or fabricates relationships. This demonstrates the need for vision encoders to be able to capture more intricate spatial details and focus less on recognition, which it already excels at but focus rather more on extracting these key spatial and visual-logical relations.



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



**Q:** Which set does the Figure belong to? Select from A, B, and C. (A) Set A (B) Set B (C) Neither set A nor set B

**CoT Response from Claude Sonnet 3.5:**

To determine which set the Figure belongs to, we need to analyze the characteristics of Set A, Set B, and the Figure:

- Set A:**
  - Contains 5 boxes
  - Each box has a mix of different sized stars
  - The number of stars in each box varies, but generally includes 4-7 stars
- Set B:**
  - Contains 5 boxes
  - Most boxes have a mix of different sized stars
  - The first box contains only one large star
  - The number of stars in each box varies, generally 3-6 stars
- The Figure:**
  - Contains exactly 4 stars
  - All stars are of similar size
  - Stars are arranged in a 2x2 grid pattern

Comparing the Figure to both sets:

- It doesn't match the typical pattern of Set A, which usually has more stars and varying sizes.
- It doesn't fit with Set B either, as it lacks the variation in star sizes seen in most boxes of Set B.
- The Figure's uniform arrangement and consistent star size is not characteristic of either set.

Therefore, the correct answer is:  
(C) Neither set A nor set B

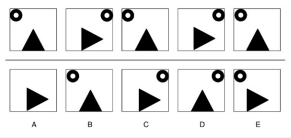
The Figure has its own unique characteristics that don't align with the patterns observed in either Set A or Set B.

**Extracted MCQ Choice:** C

**Correct MCQ Choice:** A

Figure 9: An example of Claude 3.5 Sonnet's failure to capture intricate spatial relationships is shown here. The model recognizes the differently shaped stars but misses the key relationships needed to solve the problem. It also focuses on irrelevant details, overlooking the fact that A contains an even number of stars and B an odd number, which is essential for determining the correct solution.

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



**Q:** Which set does the Figure belong to? Select from A, B, and C. (A) Set A (B) Set B (C) Neither set A nor set B

**CoT Response from Claude Sonnet 3.5:**

To determine which box comes next in the sequence, let's analyze the pattern in the given sequence:

- The sequence alternates between triangles pointing up (▲) and triangles pointing right (►).
- The circle(s) move position in each step:
  - Start in top-left corner
  - Move to top-right corner
  - Split into two circles in top-right and top-left corners**
  - Remain as two circles but move to top-right corner**
  - Return to a single circle in top-right corner**
- The last box in the given sequence has a triangle pointing up (▲) and **a single circle in the top-right corner.**

**Following this pattern, the next box should have:**

- A triangle pointing right (►)**
- Two circles in the top-right corner**

Looking at the options A through E:

A: Has the correct triangle, but no circles B: Has the wrong triangle and wrong circle position C: Has the correct triangle, but only one circle in the wrong position D: Has the wrong triangle and wrong number of circles E: Has the correct triangle (►) and the correct number and position of circles

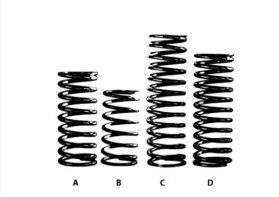
Therefore, the correct answer is E. It continues the established pattern by having a right-pointing triangle and two circles in the top-right corner.

**Extracted MCQ Choice:** E

**Correct MCQ Choice:** C

Figure 10: Another example of Claude 3.5 Sonnet’s failure to capture detailed visual logic is its hallucination of spatial facts about the relationship between the triangle and the circle. Claude incorrectly assumes that the circle is splitting in half, when in reality, this appearance results from a sequential relationship between the images.

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



**Q:** Describe this image in detail, focus on detailed patterns and spatial relationships

**CoT Response from Claude Sonnet 3.5:**

This image shows four different coil springs labeled A through D, arranged side by side from left to right. Let me describe each in detail:

**Spring A:**

- Approximately 10-11 coils
- Medium height
- Medium coil diameter
- Uniform spacing between coils

**Spring B:**

- Approximately 8-9 coils
- Shortest of the four springs
- Similar coil diameter to A
- Uniform spacing between coils

**Spring C:**

- Approximately 14-15 coils
- Tallest spring
- Similar coil diameter to A and B
- Uniform spacing between coils
- Appears slightly thinner in wire gauge

**Spring D:**

- Approximately 12-13 coils
- Second tallest spring
- Similar coil diameter to others
- Uniform spacing between coils
- Some slight irregularity at the top coil

**Spatial relationships:**

- Springs are arranged in ascending order of height from B to C
- Spring A breaks this pattern, being taller than B but shorter than C
- All springs have similar diameters
- The springs are evenly spaced horizontally
- Each spring is labeled with a letter beneath it

All springs appear to be drawn in a somewhat sketchy, hand-drawn style with some minor irregularities in the lines, giving them a more realistic rather than perfectly geometric appearance.

The overall arrangement suggests this might be a comparison or educational diagram showing different spring configurations or sizes.

Figure 11: Another example of Claude 3.5 Sonnet’s failure to capture detailed visual logic is its hallucination of spatial facts about the relationship between the triangle and the circle. Claude incorrectly assumes that the circle is splitting in half, when in reality, this appearance results from a sequential relationship between the images.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

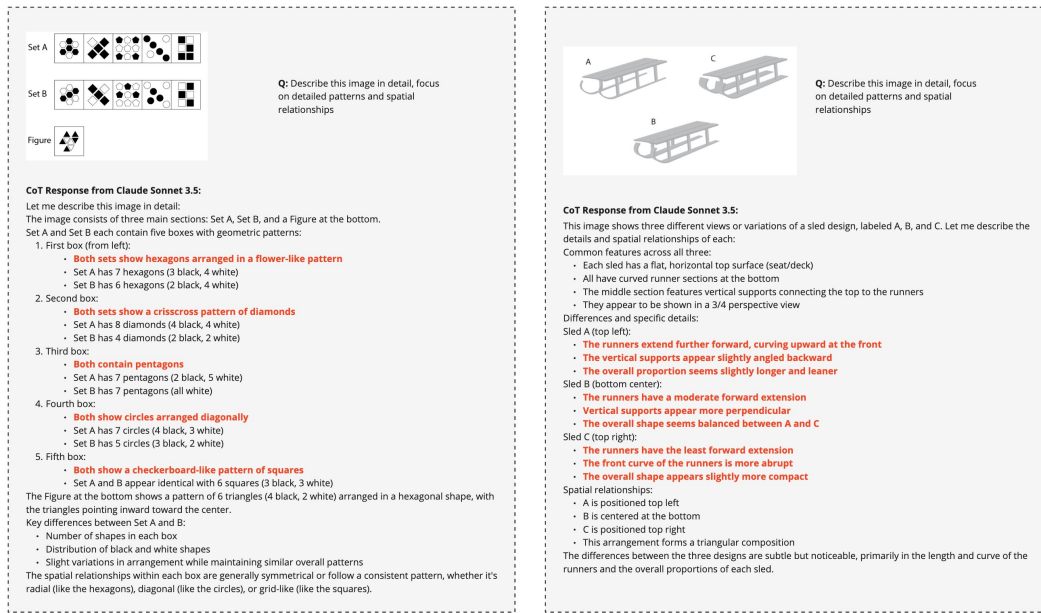


Figure 12: Another example of Claude 3.5 Sonnet’s failure to capture detailed visual logic is its hallucination of spatial facts about the relationship between the triangle and the circle. Claude incorrectly assumes that the circle is splitting in half, when in reality, this appearance results from a sequential relationship between the images.

## G MCQ VS CoT EVALUATIONS

We found that when MLLMs are prompted to produce CoT reasoning, they often provide incorrect reasoning, leading to lower scores, as open-ended evaluations account for the quality of the reasoning itself.

According to the benchmark outcome, we show that MCQ-based evaluations consistently result in higher raw scores compared to CoT-based evaluations in almost all categories for reasoning and capabilities. This mirrors human behavior, where it is often easier to make an educated guess and get the answer right without fully understanding the question or providing correct reasoning. Since MCQ evaluations may overlook cases where MLLMs guess the answer correctly without valid reasoning, we argue that CoT evaluations offer a more reliable measure of MLLM reasoning capabilities, as they assess both the answer and the reasoning behind it.

## H MODEL SIZE AND PERFORMANCE

We also observe that performance on LogicVista generally increases with model parameter sizes. As illustrated in Figure 13, there is a positive correlation between model size and average LogicVista performance. This trend suggests that larger models may possess greater capacities for learning and understanding complex relationships, allowing them to better tackle the demands of visual logical reasoning tasks. This improvement may be attributed to their ability to capture more intricate patterns and nuances in data, which enhances their overall reasoning capabilities. However, it is important to note that while larger models tend to perform better, this does not guarantee that all larger models will excel equally, as other factors such as training data quality and model architecture also play significant roles in determining performance.

## I EXAMPLES OF LOGICVISTA LOGICAL REASONING DATA

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

Table 6: Three samples requiring inductive logical reasoning skills.

(a)	
<p><b>Q:</b> Which choice (A, B, C, or D) completes the series?  <b>Answer:</b> D  <b>Reasoning:</b> In this example, there are two rules to be applied. The first is that the circle moves counter-clockwise in the hexagon. It follows that, in the following diagram, the circle will be in the upper corner of the hexagon, pointing to D as the answer. To confirm this, the second rule can be applied, according to which the position of the black triangle alternates between bottom left and top right. Thus, in the following diagram, the black triangle will need to be in the upper right corner of the hex. The answer is therefore definitely D.</p>	
<p><b>Logical Reasoning Skill:</b> Inductive  <b>Required capability:</b> Diagram</p>	
(b)	
<p><b>Q:</b> Who is the odd-one-out? Select answers from A-I.  <b>Answer:</b> G  <b>Reasoning:</b> Element G constitutes the exception and is therefore the correct answer.  <b>Logical Reasoning Skill:</b> Inductive  <b>Required capability:</b> Diagram</p>	
(c)	
<p><b>Q:</b> Two grids containing colored symbols and following a common rule are presented. In the block on the right, four additional grids are presented. The candidate must find the two grids that follow the same rule out of these four options. What options (A, B, C, or D) follow this same rule?  <b>Answer:</b> B, D  <b>Reasoning:</b> In this example, it is easy to see that the rule governing the two grids on the left is: blue triangles are present in each of the two bottom lines. This rule is followed in the two grids on the right.  <b>Logical Reasoning Skill:</b> Inductive  <b>Required capability:</b> Diagram, OCR</p>	

Table 7: Three samples requiring deductive logical reasoning skills.

1566	
1567	
1568	
1569	(a)
1570	<i>All footballers are fit and healthy.</i>
1571	
1572	<i>All famous sports players are footballers.</i>
1573	
1574	<i>Given that the above is true, which of the following is the logical deduction?</i>
1575	
1576	1. <i>All footballers are famous sports people</i>
1577	2. <i>All famous people are fit and healthy</i>
1578	3. <i>All famous sports players are fit and healthy</i>
1579	4. <i>All fit and healthy people are footballers</i>
1580	5. <i>All football players are men</i>
1581	
1582	<b>Q:</b> Which is the correct answer according to the image? Select from 1-5?
1583	<b>Answer:</b> 3
1584	<b>Reasoning:</b> Using deductive reasoning, the only logical answer is 3. To get to this answer, you
1585	need to simplify the given facts. All famous sports players are footballers, and all
1586	footballers are fit and healthy. We can not deduce that all footballers are famous
1587	sports people, as we have not got that information. We can not deduce that all
1588	famous people are fit and healthy, because the fact is about famous sports people.
1589	This is the logical answer. This information is not given; all footballers are fit and
1590	healthy but we can not logically link that all fit and healthy people are footballers.
1591	This is obviously incorrect, as gender is not mentioned at all in the question.
1592	<b>Logical Reasoning Skill:</b> Deductive
1593	<b>Required capability:</b> OCR
1594	(b)
1595	The vast majority of swallows are blue. What is the most logical conclusion?
1596	
1597	A. There is a white swallow.
1598	B. Not everything that is blue is a swallow.
1599	C. There is a blue swallow.
1600	D. None of the answers are satisfactory.
1601	
1602	<b>Q:</b> What is the correct answer to the question in the image? Select from A-D?
1603	<b>Answer:</b> C
1604	<b>Reasoning:</b> The vast majority of swallows are blue so the answer must be C: there is a blue
1605	swallow.
1606	<b>Logical Reasoning Skill:</b> Deductive
1607	<b>Required capability:</b> OCR
1608	(c)
1609	<p>The people determine what is produced.  The government is made up of the people.  Production is determined by the free-market.  The free-market is made up of production.  Government is determined by the free-market.</p>
1610	
1611	
1612	<b>Q:</b> What is produced is determined by the people. Select from A, B and C. (A) True
1613	(B)False (C)Insufficient Information?
1614	<b>Answer:</b> A
1615	<b>Reasoning:</b> Line 1 states that the people determine what is produced. Line 2 states that the
1616	government is made up of the people. Therefore, the people determine what is
1617	produced. This is a syllogism. Thus, this statement is true.
1618	<b>Logical Reasoning Skill:</b> Deductive
1619	<b>Required capability:</b> OCR

Table 8: Three samples requiring numerical logical reasoning skills.

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

(a)

Share Price Index				
Company	Today's Price (€)	Change from previous day (%)	Past 12 months Max price (€)	Past 12 months Min price (€)
Huver Co.	1,150	1.10	1,360	860
Drebs Ltd	18	0.50	22	11
Fevs Plc	1,586	-9.00	1,955	1,242
Fauvers	507	-1.00	724	464
Steapars	2,537	1.00	2,630	2,216

Dividend Index					
Dividend paid per share (€)	Huver Co.	Drebs Ltd	Fevs Plc	Fauvers	Steapars
Interim Dividend	0.83	0.44	0.34	0.09	0.48
Final Dividend	1.75	1.12	1.25	0.32	0.96

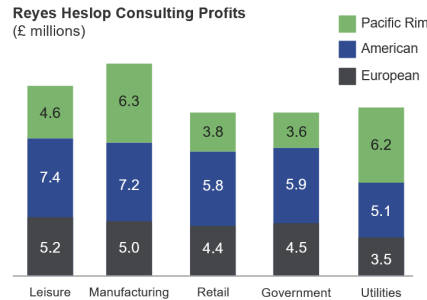
**Note:** the total annual dividend paid per share is the sum of the interim dividend and the final dividend.

**Q:** Which share had the largest difference between highest and lowest price over the last 12 months? Select from A, B, C, D and E. (A) Huver Co. (B) Drebs Ltd (C) Fevs Plc (D) Fauvers (E) Steapars

**Answer:** C  
**Reasoning:** Step 1- Calculate the difference between the maximum and the minimum prices. Huver Co. = 1,360 - 860 = 500 Drebs Ltd = 22 - 11 = 11 Fevs Plc = 1,955 - 1,242 = 713 Fauvers = 724 - 464 = 260 Steapars = 2,630 - 2,216 = 414. Tip: Notice the wording of the question is asking for the share with the largest absolute change in price, NOT the largest percentage change, which would have been Drebs Ltd. If the question had wanted the percentage change it would have used the word percentage. Thus the correct answer is (C) Fevs Plc

**Logical Reasoning Skill:** Numerical  
**Required capability:** OCR

(b)

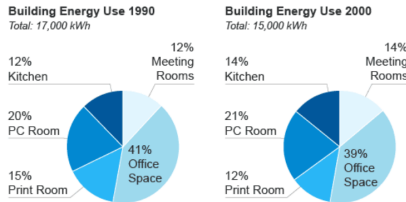


**Q:** Reyes Heslop had a target for Leisure profits to be a quarter of their total profits. Assuming profits in other areas remain the same, by how much did the Leisure profits miss this target? Select from A, B, C, D and E. (A) 31.8 million (B) 32.4 million (C) 32.7 million (D) 33.2 million (E) 33.4 million

**Answer:** D  
**Reasoning:** Step 1- Calculate the total Reyes Heslop profits across all areas other than Leisure. (6.3 + 7.2 + 5.0) + (3.8 + 5.8 + 4.4) + (3.6 + 5.9 + 4.5) + (6.2 + 5.1 + 3.5) = 61.3 million. Step 2- This needs to be 1/4 of all profits for the condition to be met. Therefore all profits, across all sectors, would be 61.3 / 75% = 81.7333 million. Step 3- Now we look at the difference between actual and target Leisure profits. Actual = (4.6 + 7.4 + 5.2) = 17.2 Target = (81.7333 - 61.3) = 20.4333 Shortfall = 3.2333 (millions) Thus the correct answer is (D) 33.2 million

**Logical Reasoning Skill:** Numerical  
**Required capability:** Diagram, OCR

(c)



**Q:** Which space experienced the smallest reduction in kWh used between 1990 and 2000? Select from A, B, C, and D. (A) Office Space (B) Print Room (C) Meeting Rooms (D) PC Room

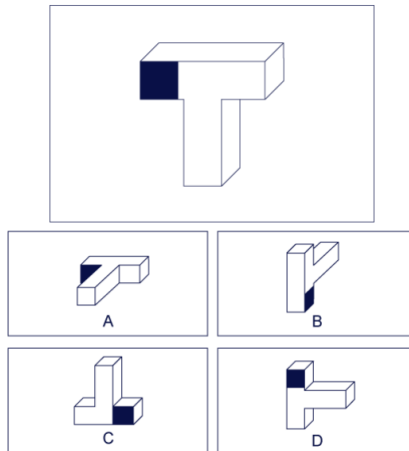
**Answer:** D  
**Reasoning:** Step 1- Calculate the value of kWh for 1990 and 2000 for each of the rooms. Room 1990 per kWh 2000 per kWh Meeting Rooms 2.04 2.10 Office Space 6.97 5.85 Print Room 2.55 1.80 PC Room 3.40 3.15 Kitchen 2.04 2.10 Step 2- Subtract the kWh for 2000 from that of 1990 for each of the rooms. Room change (1990 - 2000) kWh Meeting Rooms -0.06 Office Space 1.12 Print Room 0.75 PC Room 0.25 Kitchen -0.06 Step 3- Look for the smallest positive value. Negative values represent an increase between 1990 and 2000. Tip- You only need to perform 4 calculations, as two of the rooms have the same values. Thus, the correct answer is (D) PC Room.

**Logical Reasoning Skill:** Deductive  
**Required capability:** Diagram, OCR

Table 9: Three samples requiring spatial logical reasoning skills.

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

(a)

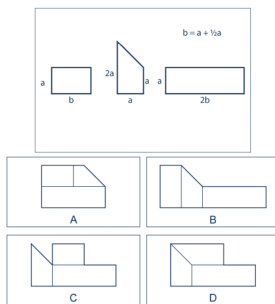


**Q:** Which figure is a rotation of the object? Select from A, B, C, and D. (A) (B) (C) (D)

**Answer:** B  
**Reasoning:** The answer is B.

**Logical Reasoning Skill:** Spatial  
**Required capability:** Diagram

(b)

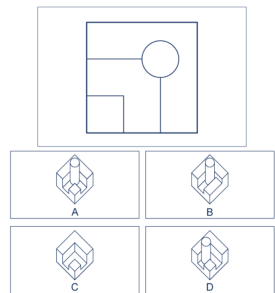


**Q:** Which figure can be formed with the given piece? Select from A, B, C, and D. (A) (B) (C) (D)

**Answer:** C  
**Reasoning:** The answer is C.

**Logical Reasoning Skill:** Spatial  
**Required capability:** Diagram

(c)



**Q:** To which object does the given top view correspond? Select from A, B, C, and D. (A) (B) (C) (D)



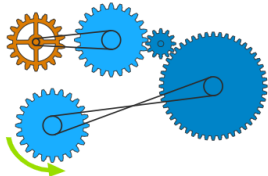
**Answer:** A  
**Reasoning:** The answer is A.

**Logical Reasoning Skill:** Spatial  
**Required capability:** Diagram



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

Table 10: Three samples requiring mechanical logical reasoning skills.

(a)	
<p><b>Q:</b> A non-pressurised cylindrical metal tank filled with air is submerged underwater. As the air escapes, the tank gradually moves deeper underwater. Which statement provides the best reason for this motion? Select from A, B, C, D, and E. (A) The bubbles provide a downward thrust on the tank (B) The metal increases in density so it gets heavier (C) The bubbles lower the density of the water which lowers its buoyancy (D) Water replaces the air in the tank which makes it heavier (E) Impossible to tell</p> <p><b>Answer:</b> D</p> <p><b>Reasoning:</b> As air escapes the available space is quickly replaced with water, so the tank's density becomes the same as that of the water and with the added weight and density of the tank itself continues to sink.</p> <p><b>Logical Reasoning Skill:</b> Mechanical</p> <p><b>Required capability:</b> Diagram</p>	
(b)	
<p><b>Q:</b> It is a cold winter outside and a well-insulated house has its heater turned on. The front door is opened and cold air rushes in. If the wind speed outside is very low, how would the cold air enter the house? Select from A, B, C, D, and E. (A) Scenario A, the cold air will flow towards the floor (B) Scenario B, the cold air will flow towards the ceiling (C) A combination of A and B (D) The cold air will not enter the house (E) Impossible to tell</p> <p><b>Answer:</b> A</p> <p><b>Reasoning:</b> Cold air sinks, whereas hot air rises. The house and the air inside it are warmer than the outside air temperature, so if these two systems (house and outside) were to be suddenly connected (door opening) the cold air would sink and the hot air would sit above the cold air until the heat transferred between the two.</p> <p><b>Logical Reasoning Skill:</b> Mechanical</p> <p><b>Required capability:</b> Diagram</p>	
(c)	
<p><b>Q:</b> In which direction does the orange gear rotate? Select from A, B, and C. (A) Clockwise (B) Counterclockwise (C) No rotation</p> <p><b>Answer:</b> A</p> <p><b>Reasoning:</b> The correct answer is clockwise.</p> <p><b>Logical Reasoning Skill:</b> Mechanical</p> <p><b>Required capability:</b> Diagram</p>	

1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835

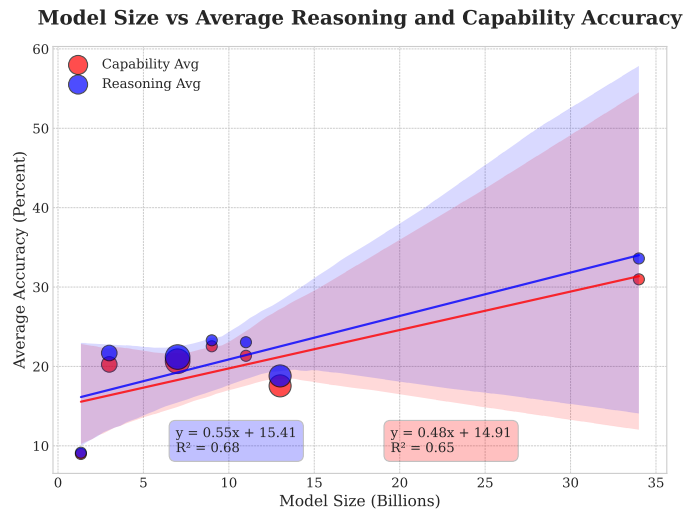
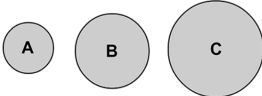
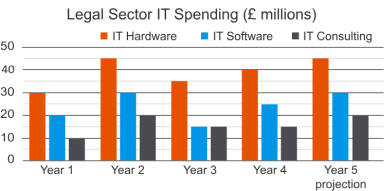
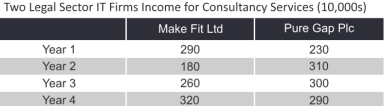


Figure 13: Correlation between Model Size and Average Accuracy. The scatter plot employs varying dot sizes to indicate the number of models with identical model sizes, illustrating the distribution density.

## J EXAMPLES OF DIFFERENT BROAD LOGICVISTA CAPABILITIES DATA

Table 11: Three samples of diagram, OCR, and mixed LogicVista data

1836	
1837	
1838	
1839	
1840	
1841	(a)
1842	
1843	
1844	
1845	<b>Q:</b> Which ball is the heaviest? Select from A, B, C, and D. (A) A (B) B (C) C (D) CAN
1846	NOT SAY
1847	<b>Answer:</b> D
1848	<b>Reasoning:</b> The correct answer is D.
1849	<b>Logical Reasoning Skill:</b> Mechanical
1850	<b>Required capability:</b> Diagram
1851	(b)
1852	Which of these objects will not float on water?
1853	
1854	
1855	<b>Q:</b> Select from A, B, C, and D. (A) banana (B) scissors (C) empty plastic soda bottle
1856	(D) wooden pencil
1857	<b>Answer:</b> B
1858	<b>Reasoning:</b> The correct answer is B because scissors have metal and are most likely to sink.
1859	<b>Logical Reasoning Skill:</b> Deductive
1860	<b>Required capability:</b> OCR
1861	(c)
1862	
1863	
1864	
1865	
1866	
1867	
1868	
1869	
1870	
1871	<b>Q:</b> Which of the following statements is false regarding legal sector spending between
1872	Year 4 and projected Year 5? Select from A, B, C, D and E. (A) IT consulting will
1873	increase by 35million. (B) IT consulting will match that of year 2. (C) IT software
1874	will exceed IT consulting. (D) Spending on IT hardware will decline. (E) None of
1875	these.
1876	<b>Answer:</b> D
1877	<b>Reasoning:</b> Step 1- Check in turn whether each statement is true or false: a) The projected
1878	spend on IT consulting is projected to increase by 35 million. Option A is true. b)
1879	The projected spend on IT consulting is 320 million, which matches year 2. Option
1880	B is true. c) The projected spend on IT software is 330 million and for IT consulting
1881	it is 320 million. Option C is true. d) There are increases projected for IT hardware,
1882	for IT software and for consulting, therefore spending on IT hardware will decline is
1883	not true. The option for D is false. e) We see that option D is false, so E cannot
1884	be the correct answer. Thus the correct answer is (D) Spending on IT hardware,
1885	software and consulting is projected to decline.
1886	<b>Logical Reasoning Skill:</b> Numerical
1887	<b>Required capability:</b> Diagram, OCR
1888	
1889	

## K DATA LEAKAGE CONCERNS OF EXISTING BENCHMARKS

As shown in 14, sourcing data from gated sources allows LogicVista to greatly minimize the risk of data leakage. In an experiment, we randomly sampled 50 images from datasets such as MM-vet,

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

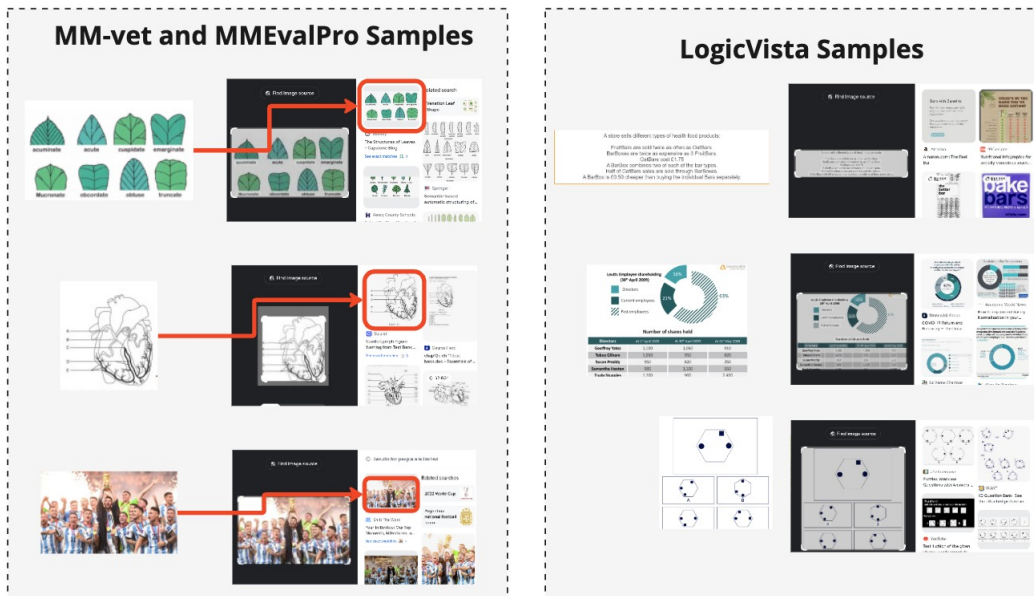


Figure 14: LogicVista mitigates potential data leakage by sourcing from gated private datasets (with permission). This approach ensures a fair comparison by isolating MLLM reasoning abilities, preventing any overlap with information that may have been included in their training data.

MMEvalPro, RAVEN, and MathVista Yu et al. (2023); Huang et al. (2024); Zhang et al. (2019); Lu et al. (2024), and used Google’s reverse image search. We found that all samples from existing benchmarks were publicly available online, whereas nearly all of LogicVista’s samples were inaccessible, either behind paywalls or requiring registration. Since most of LogicVista’s data is not publicly available, it is much more difficult to scrape for training MLLM models. This restricted access reduces the chances of LogicVista’s samples being included in training datasets, unlike in open benchmarks.

## L CROWDSOURCING ANNOTATION TOOL

To scale LogicVista for the future, we have released an annotation tool similar to the one used in our annotation process. This tool facilitates robust annotations by incorporating rounds of peer review before finalizing entries in LogicVista. Additionally, it is web-based, allowing the community to contribute to LogicVista from anywhere. We hope this will enable LogicVista to grow and increase its sample size significantly.

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

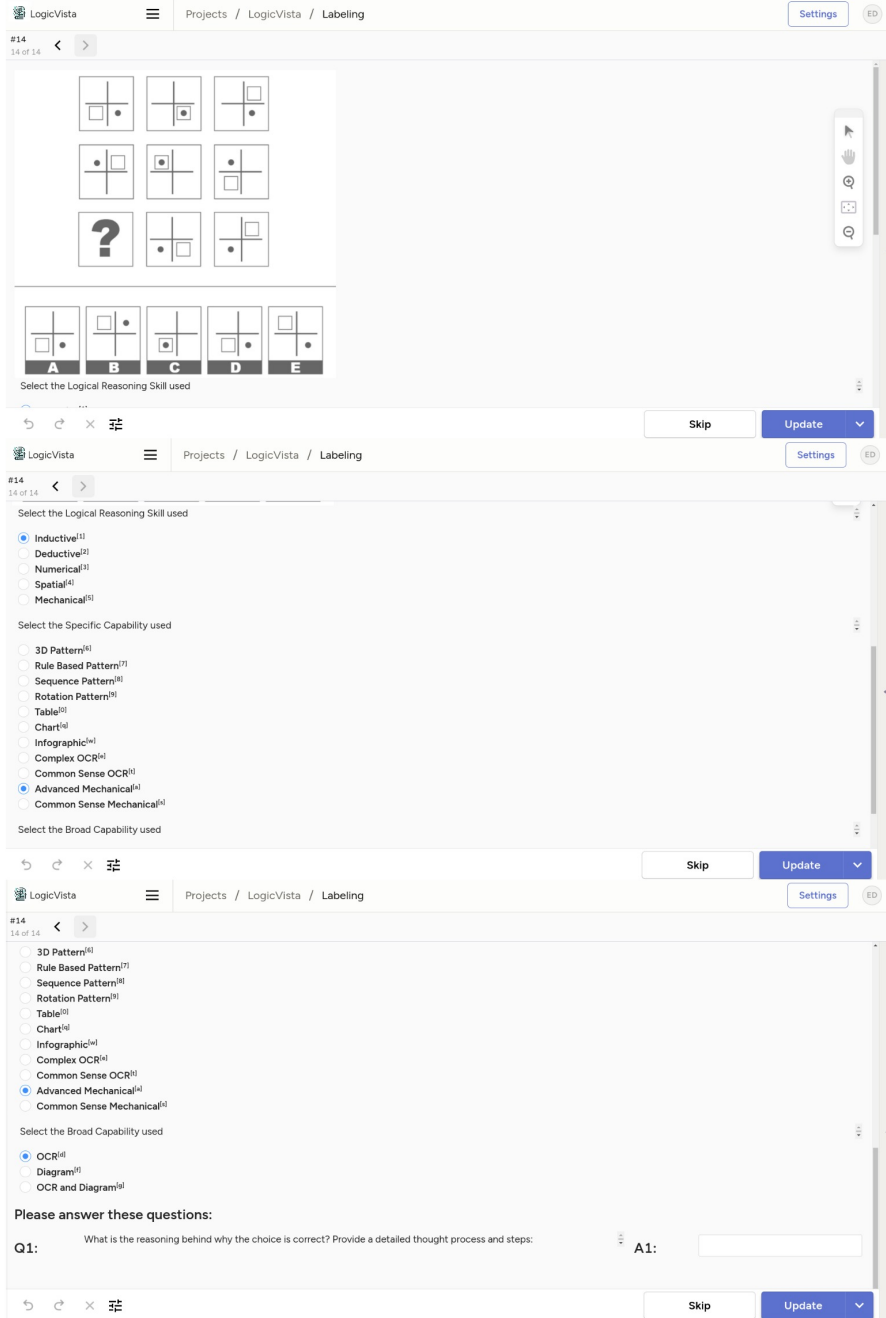


Figure 15: Example of the annotation process using our tool, enabling the community to contribute to scaling LogicVista effectively.

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

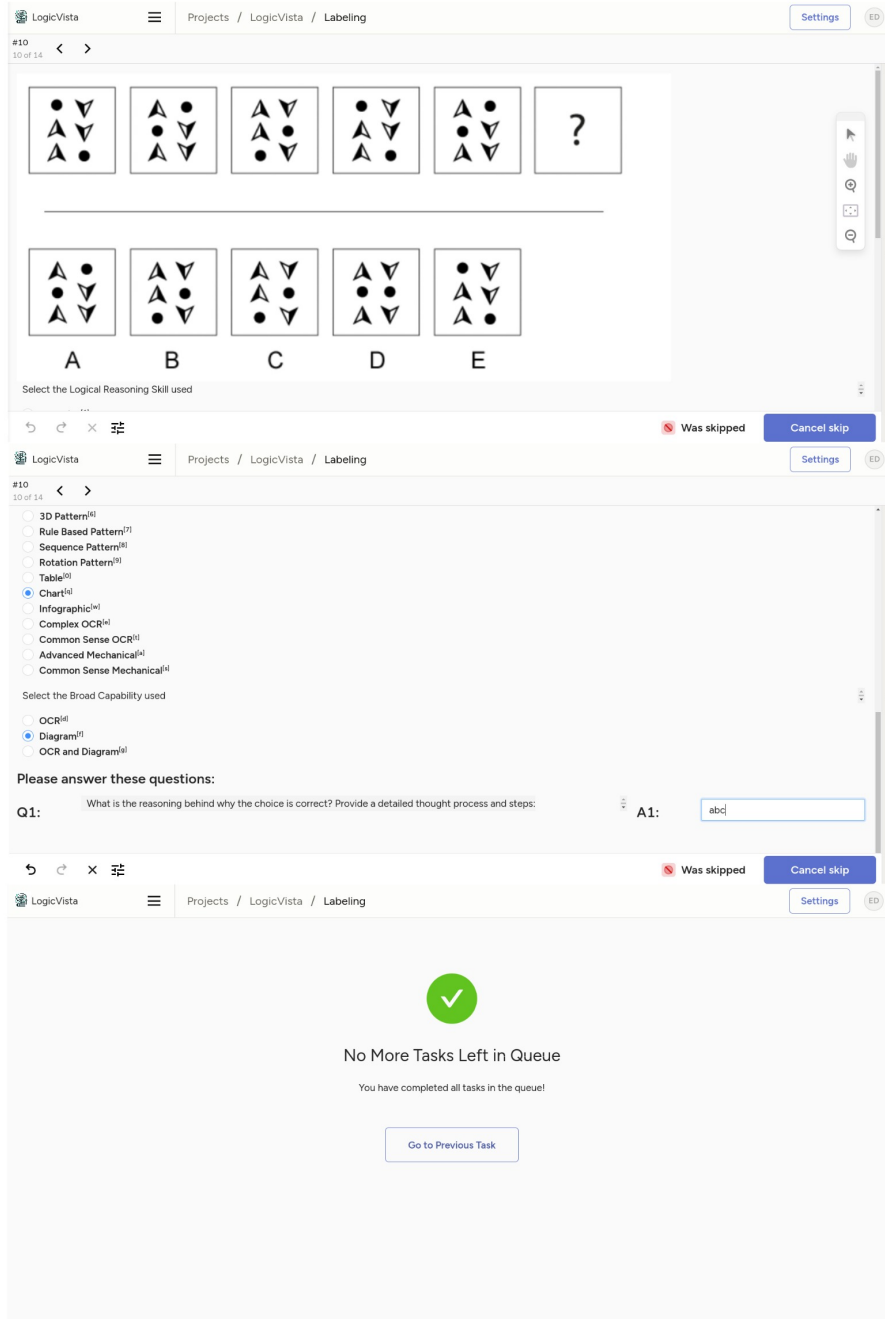


Figure 16: Additional example of annotation process using our crowdsourcing tool