
VisuLogic: A Benchmark for Evaluating Visual Reasoning in Multi-modal Large Language Models

Weiyi Xu^{1,3*}, Jiahao Wang^{2,3*}, Weiyun Wang^{3†}, Zhe Chen^{3†}, Wengang Zhou¹, Aijun Yang², Lewei Lu⁴, Houqiang Li¹, Xiaohua Wang², Xizhou Zhu³, Wenhai Wang³, Jifeng Dai^{5,3✉}, Jinguo Zhu^{3✉}

¹University of Science and Technology of China, ²Xi'an Jiaotong University,

³Shanghai Artificial Intelligence Laboratory, ⁴SenseTime Research, ⁵Tsinghua University

ustcxwy0271@mail.ustc.edu.cn, wjhwds@stu.xjtu.edu.cn, lechatelia@gmail.com

Abstract

Visual reasoning is a core component of human intelligence and a critical capability for advanced multimodal models. Yet current reasoning evaluations of multi-modal large language models (MLLMs) often rely on text descriptions and allow language-based reasoning shortcuts, failing to measure genuine vision-centric reasoning. To address this, we introduce VisuLogic: a benchmark of 1,000 human-verified problems across six categories (e.g., quantitative shifts, spatial relations, attribute comparisons). These various types of questions can be evaluated to assess the visual reasoning capabilities of MLLMs from multiple perspectives. We evaluate leading MLLMs on this benchmark and analyze their results to identify common failure modes. Most models score below 30% accuracy—only slightly above the 25% random baseline and far below the 51.4% achieved by humans—revealing significant gaps in visual reasoning. Furthermore, we provide an supplementary training dataset and a reinforcement-learning baseline to support further progress. Code, data, and baselines are available at <https://visulogic-benchmark.github.io/VisuLogic>.

1 Introduction

Reasoning, as fundamental component of human intelligence, has become a critical criterion in evaluating progress toward Artificial General Intelligence (AGI) [25, 74]. Recent advancements in Large Language Models (LLMs) have demonstrated substantial improvements in reasoning capabilities across complex domains such as mathematics [61, 82, 81, 58], logical reasoning [68, 79, 22, 47] and coding [2, 34, 42, 31]. Techniques like Chain-of-Thought (CoT) [75] prompting and test-time compute scaling (e.g., OpenAI o1 [33] and Deepseek-R1 [17]) have significantly enhanced the reasoning performance of LLMs [17, 25, 74]. Along with the rapid development of language reasoning research for LLMs, considerable progress [84, 61, 58, 10, 50, 72, 51, 63, 73] has been made in improving multimodal reasoning capability of Multimodal Large Language Models (MLLMs).

These methods, which often incorporate reinforcement learning techniques [10, 50, 61] to enhance the reasoning capabilities of MLLMs, have achieved some early successes [84, 61, 58, 10, 50, 51, 63]. However, they typically rely on existing multi-modal benchmarks that struggle to accurately capture a model’s core visual reasoning ability. For example, VLM-R1 [63] assesses “visual reasoning” with referring expression comprehension tasks [88, 55, 38], yet these tasks primarily focus on object localization, demanding only basic perceptual skills rather than more advanced visual cognitive processes. Meanwhile, several works [58, 61, 84] adopt mathematical problem-solving benchmarks that include diagrams—such as MathVista [52], MathVerse [91], and MATH-Vision [69]—to evaluate visual reasoning. In practice, however, as [91] observes, many MLLMs translate these visual clues

*equal contribution; † interns at OpenGVLab, Shanghai AI Laboratory; ✉ corresponding author.

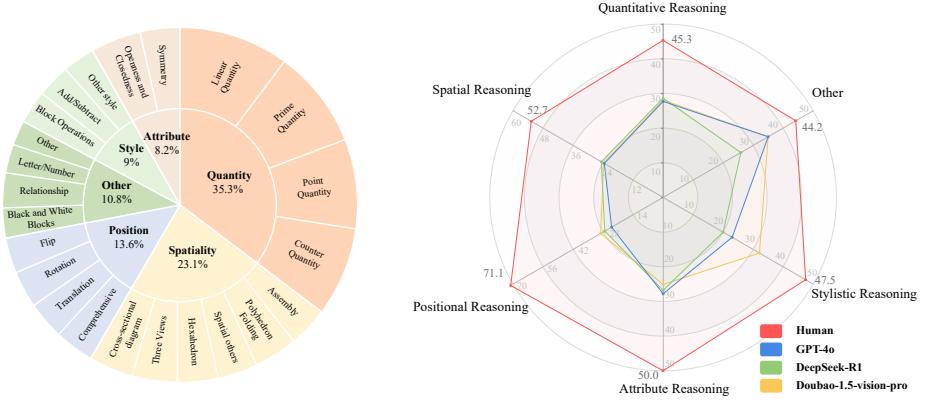


Figure 1: Composition of the VisuLogic benchmark and performance of representative MLLMs.
The left figure shows the distribution of the 6 categories and their subcategories in VisuLogic. The right figure shows accuracies (%) achieved by MLLMs and by human on each category of VisuLogic.

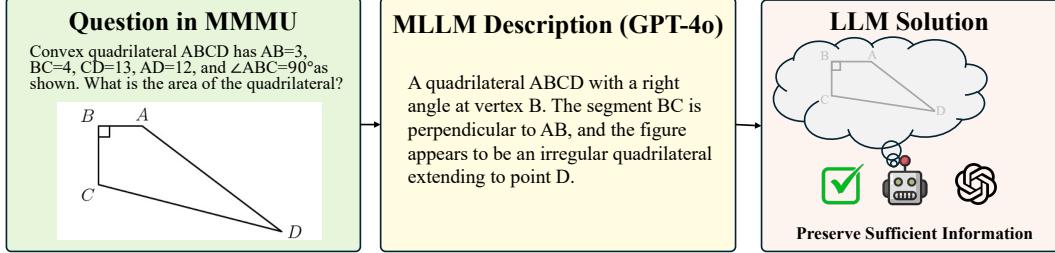
into textual descriptions and then rely on standard language reasoning. This approach can incorrectly attribute language-driven results to visual reasoning, resulting in a misleading assessment of the model’s visual reasoning capabilities [91, 29]. Consequently, designing new benchmarks that explicitly focus on vision-centric reasoning—rather than conflating it with text-based reasoning—remains critical for advancing MLLMs’ visual reasoning capacities.

To address this limitation, we propose VisuLogic, a novel benchmark specifically designed to evaluate visual reasoning abilities in multimodal models without mixing them with purely text-based reasoning. VisuLogic comprises carefully constructed tasks that span multiple reasoning categories, as shown in Figure 1. As shown in Figure 5, these tasks are classified into six key types, such as Quantitative Reasoning, which requires understanding and deducing shifts in the quantity of certain elements within an image (see Figure 3). In contrast to existing benchmarks, as demonstrated in Figure 2, state-of-the-art (SOTA) MLLMs often omit crucial visual details when describing VisuLogic problems, making it difficult for them to rely solely on a text-based inference shortcut. Indeed, even humans would find it challenging to capture every essential visual cue in a single description, so effectively tackling VisuLogic demands more robust, vision-centric reasoning. By reducing reliance on textual inference shortcuts, VisuLogic thus provides a stringent evaluation of MLLMs’ genuine visual reasoning capabilities.

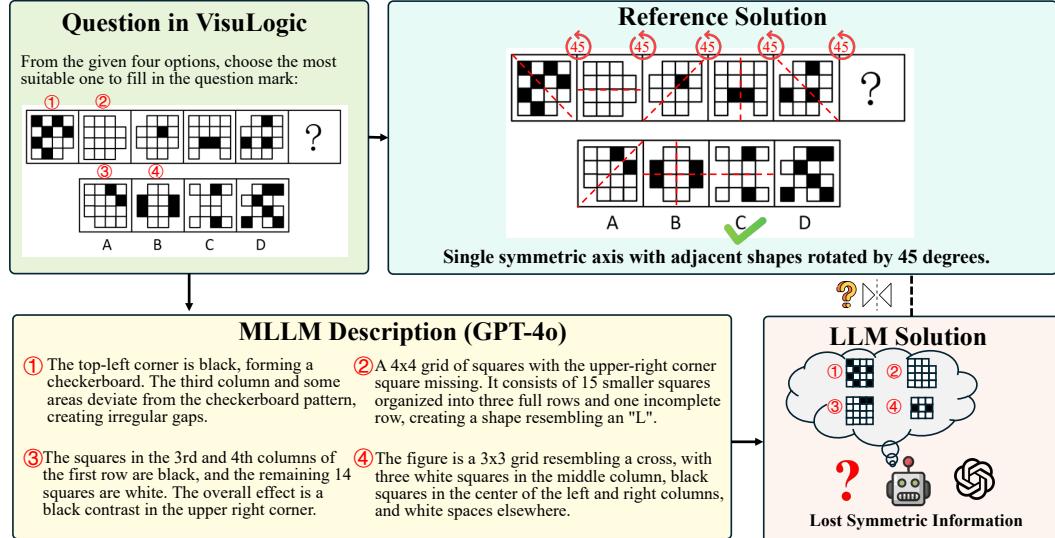
We conducted a comprehensive evaluation and systematic analysis to assess current models’ visual reasoning capabilities. When leading text-only LLMs were supplied with detailed descriptions in place of raw images, their accuracy—Doubao-1.5-Pro (26.6%), Claude-3.7-Sonnet (25.9%) and Qwen2.5-72B-Instruct [83] (28.0%)—barely exceeded the random-chance baseline of 24.9%. This clearly demonstrates that textual reasoning alone are insufficient for solving our VisuLogic tasks. Even state-of-the-art multimodal LLMs (MLLMs)—including GPT-4o [32], Doubao-1.5-Vision-Pro, Gemini-2.0-Pro-Exp [64] and InternVL3-78B [94]—achieved only 26.3%, 28.1%, 28.0% and 27.7%, respectively, whereas human participants reached 51.4%. The substantial gap between these results and human performance underscores the challenge of robust visual reasoning in current MLLMs. Furthermore, we applied a simple reinforcement-learning (RL) fine-tuning step on our supplementary training dataset: this boosted the baseline model’s accuracy from 25.5% to 31.1%, outperforming both open-source and closed-source counterparts. These findings illustrate the promise of the RL technique for strengthening MLLMs’ visual reasoning capabilities.

In summary, our contributions are as follows:

- We propose a challenging visual reasoning benchmark that is inherently difficult to articulate using language, providing a more rigorous evaluation of the visual reasoning capabilities of MLLMs.
- We conduct comprehensive experiments to evaluate and analyze the benchmark, including extensive evaluations and comparative studies of various MLLMs under different setting.



(a) Pipeline of “MLLM description→LLM” for Question in MMMU [89]. It is trivial that SOTA MLLMs extract key visual details, thereby enabling the LLM to answer questions solely based on language reasoning.



(b) Pipeline of “MLLM description→LLM” for Question in VisuLogic. Even SOTA MLLMs struggle to describe images precisely, leading to ambiguous interpretations.

Figure 2: **Comparison of the “MLLM description→LLM” pipeline on two benchmarks.** In MMMU, detailed descriptions lead to correct solutions, while in VisuLogic, critical visual cues (*e.g.*, symmetry, rotation) can be easily lost, causing the LLM to misinterpret the pattern. This highlights that textual reasoning alone is insufficient, underscoring the benchmark’s demand for robust and in-depth visual reasoning.

- We identify the RL technique as a promising direction for improving the visual reasoning capabilities of MLLMs. Furthermore, we release both the training code and data to facilitate future research.

2 Related Work

Multi-modal Large Language Models. Recent years have witnessed substantial advancements in Multi-modal Large Language Models (MLLMs). Early works like BLIP [41, 40] and Flamingo [4] introduce lightweight parameters between vision transformer [20] (ViT) and LLMs, laying the groundwork for multimodal perception. Subsequent efforts, such as LLaVA [45] and MiniGPT-4 [93], integrate instruction tuning, further enhancing the performance of MLLMs. Proprietary models like GPT-4o [32] and Gemini-Pro [64] have advanced MLLM performance on complex multimodal tasks, while open-source models such as Qwen-VL series [6, 70, 7] and InternVL series [14, 15, 23, 13, 94] achieve competitive results through optimized architectural design, dataset expansion and training paradigm improvements. Meanwhile, some related studies further advance the ability of large models by incorporating new modalities (*e.g.*, audio [21, 18, 77], point clouds [26, 8], video [92, 11]) and by supporting more tasks (*e.g.*, grounding [80, 71], computer usage [60, 5]). Notably, limited research attempts to enhance the reasoning capabilities of MLLMs. Some pioneering works, such as R1-Onenvision [84], LMM-R1 [61], MM-EUREKA [58], R1-V [10], Visual-rft [50], Visualprm [72],

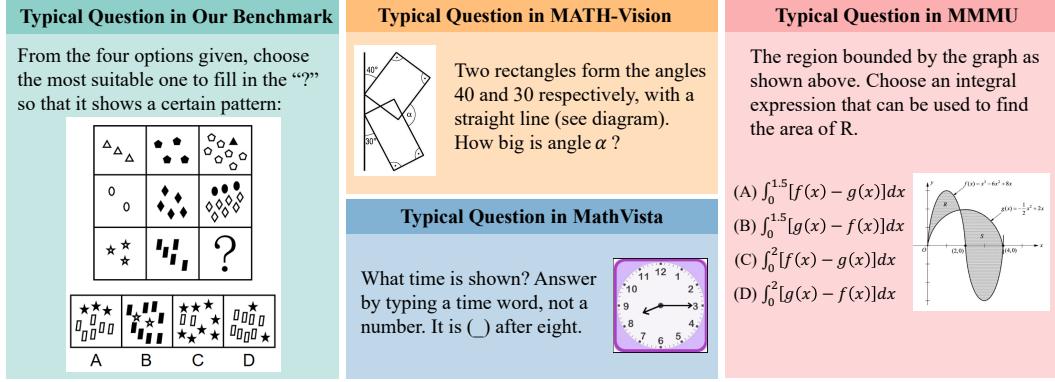


Figure 3: **Comparison of questions from different Benchmarks.** Compared to MathVista [52], MATH-Vision [69], and MMMU [89], VisuLogic focuses more explicitly on assessing pure visual reasoning capabilities.

OThink-MR1 [51], VLM-R1 [63], and Open-r1-Video [73] have explored the visual reasoning capabilities of MLLMs through Reinforcement Learning (RL), but they are still in the nascent stage.

Multimodal Benchmarks. With the development of MLLMs, multimodal benchmarks have also evolved significantly [43]. Early benchmarks primarily address visual perception tasks through simple tasks like visual question answering (VQA) [12, 44, 35, 78], image captioning [59, 19, 36] and referring expression comprehension [88, 55]. Subsequent works expand the capability coverage of benchmarks into more specialized domains: OCRBench [49], Chartqa [56] and DocVQA [57] assess textual content extraction; AgentBench [48] and ToolEyes [86] test tool usage capabilities; and egocentric perception benchmarks [54, 16] quantify first-person scene interpretation. Despite the progress, they ignore the evaluation of visual reasoning abilities [90, 89]. Recently, some benchmarks have made explorations in examining MLLMs’ visual reasoning abilities, but methodological deficiencies still cause limitations to assess the intrinsic visual reasoning capabilities [29, 3, 76]. InfiIMM-Eval [28] test reasoning abilities around daily life, lacking deep-level reasoning scenarios. MMMU [89] and Emma [29] provide benchmarks demanding advanced reasoning abilities in fields such as chemistry and physics, but they ignore questions around the images’ fundamental visual components (e.g., shapes, elements). While mathematical benchmarks [69, 52, 30, 62, 91, 27] evaluate mathematical reasoning with geometric and diagram problems included, they focus on math capabilities but disregard logical analysis about the vision information. LogicVista [76] provides a multimodal logical reasoning benchmark, its visual questions lack analytical depth—dominated by single-hop, superficial queries in limited data scope. Unlike previous works, we introduce a challenging benchmark focused specifically on the domain of visual logical reasoning.

3 VisuLogic

In this section, we first describe the VisuLogic data-curation pipeline, which comprises three key stages: data collection, quality control, and the detailed taxonomy. We then report the benchmark’s construction statistics, including total size, answer-option distributions, and category-level proportions. Finally, we introduce a supplementary training dataset—consisting of questions analogous to those in VisuLogic—designed to bolster future research and facilitate community engagement.

3.1 Data Curation Pipeline

Data Collection. We construct the VisuLogic dataset by sourcing all questions from publicly available online resources in compliance with relevant licenses and regulations. As shown in Figure 4, our automated data processing pipeline comprises three stages: 1) **Fetching**: We employ Playwright¹ to systematically scrape raw web content, supplemented by custom parsing scripts that extract question–answer pairs. 2) **Cleaning**: We remove noise, irrelevant content, and extraneous HTML markup (e.g., `<div>`) to ensure the integrity of the textual data. 3) **Structuring**: We standardize

¹<https://github.com/microsoft/playwright>

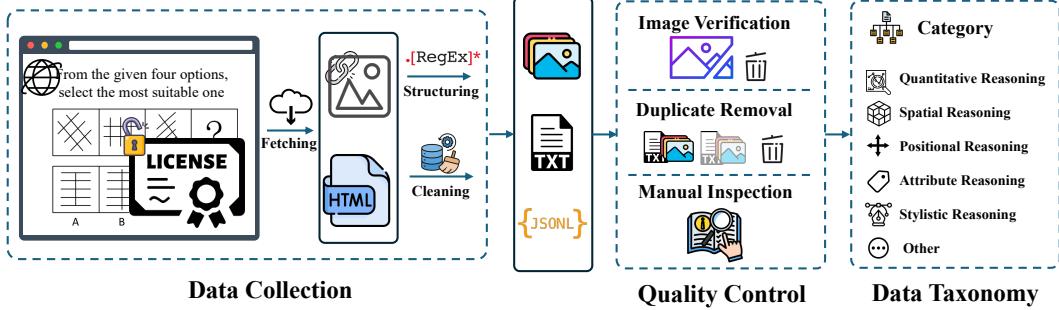


Figure 4: **Data Curation Pipeline of VisuLogic.** The pipeline includes Data Collection, Quality Control and Data Taxonomy.

the cleaned text and images by structuring all information (such as question text, category labels, and metadata) in JSON Lines (JSONL) format. To maintain visual consistency, we constrain each image’s shorter side to between 200 and 700 pixels—preserving the original aspect ratio to avoid significant size discrepancies.

Quality Control. To ensure the reliability of the benchmark dataset, we employ a three-stage data validation procedure: 1) **Image Verification:** Each image referenced in the questions is checked for existence and correct formatting; any items that fail to meet these criteria are removed following human review. 2) **Duplicate Removal:** We eliminate redundant entries at both the text and image levels by (i) detecting lexical overlap among text segments and (ii) applying perceptual hashing (pHash) to identify visually similar images. 3) **Manual Checking:** After automated filtering, we perform a thorough human-led review of every remaining entry to confirm its validity and ensure dataset reliability.

Data Taxonomy. We categorize all collected data into a taxonomy of six primary classes based on expert human annotation of the reasoning skills each question requires. Annotators first tag questions according to the targeted reasoning competency; these annotated tags are then analyzed and merged into five primary categories. A subsequent human review ensures that every question is accurately classified, with any ambiguous instances consolidated under the “Other” category. Specifically, we define each category as follows. **Quantitative Reasoning** focuses on changes in the number or count of graphical elements (for example, points, lines and angles) and on arithmetic relationships among shapes. **Spatial Reasoning** requires mentally reconstructing three-dimensional shapes from two-dimensional figures, folding or unfolding surfaces, and integrating three-dimensional structures. **Positional Reasoning** examines transformations such as translation, rotation and reflection of objects while preserving their fundamental elements. **Attribute Reasoning** involves intrinsic properties of shapes, including symmetry (axial or central), curvature and measures of openness or closedness. **Stylistic Reasoning** entails alterations in stylistic features such as overlay, subtraction and assessments of shape similarity or difference. **Other** encompasses questions that fall outside the preceding categories, including those involving letters, alphanumeric symbols or other specialized characters.

3.2 Dataset Statistics

Following data curation and validation, VisuLogic comprises 1,000 single-choice questions. Figure 1 (left) illustrates the category distribution: Quantitative Reasoning (35.3%), Spatial Reasoning (23.1%), Positional Reasoning (13.6%), Attribute Reasoning (8.2%), Stylistic Reasoning (9.0%), and Other (10.8%). Correct answer options are evenly balanced, with the proportions distributed as follows: A (23.1%), B (26.7%), C (25.2%), and D (25.0%).

3.3 Supplementary Training Dataset

To facilitate further investigation of visual reasoning, we provide an auxiliary training set of 4,296 question–answer pairs drawn from the same domains and subjected to identical validation procedures to prevent overlap with the benchmark. The training split mirrors the primary taxonomy, with category proportions of Quantitative Reasoning (30.7%), Spatial Reasoning (25.5%), Positional Reasoning (13.0%), Attribute Reasoning (8.8%), Stylistic Reasoning (9.9%), and Other (12.1%).

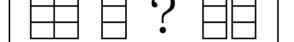
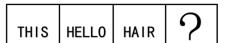
Quantitative Reasoning	Spatial Reasoning	Positional Reasoning
Choose the most suitable option to fill in the question mark, so that a certain pattern is presented.  A B C D	On the left is a net of the outer surface of a cube. Which option on the right can be formed by folding it?  A B C D	The option that matches the given pattern is ().  A B C D
Attribute Reasoning	Stylistic Reasoning	Other
Select the most appropriate one to fill in the question mark to maintain certain pattern.  A B C D	Choose the most suitable option to fill in the question mark, so that it presents a certain regularity.  A B C D	Choose the most suitable one to fill in the question mark, so that it presents a certain regularity.  A B C D

Figure 5: **Problem examples of different categories in our VisuLogic Benchmark.** VisuLogic contains 6 categories of questions, which require models’ abilities in visual logic reasoning.

4 Experiments

In this section, we present a comprehensive evaluation of the VisuLogic benchmark. We first describe the experimental setup in Section 4.1, followed by overall performance results in Section 4.2. We then analyze systematic errors in Section 4.3 and provide qualitative insights in Section 4.4.

4.1 Experiment Setup

References Performance. To contextualize model performance, we establish two reference points: 1) **Human Performance:** We invited 100 graduate students in science and engineering to solve 10 randomly sampled VisuLogic questions each, allowing 2–5 minutes per question. The aggregate accuracy over all participants constitutes the human benchmark. 2) **Random Selection:** We simulate random guessing by sampling answers uniformly over 10 independent runs and report the average accuracy as the random baseline.

Evaluated Models. We evaluate a total of 28 models on VisuLogic, comprising 8 large language models (LLMs) and 20 multimodal large language models (MLLMs). For open-source LLMs, we test *Deepseek-R1* [17], *Qwen2.5-72B-Instruct* [83] and *Qwen-QwQ* [67], and for close-source LLMs we evaluate *GPT-4* [1], *o3-mini*, *Gemini-2.0-Flash-Thinking* [64], *Claude-3.7-Sonnet* and *Doubaio-1.5-Pro-32k*. Open-source MLLMs include *Qwen2.5-VL-7B-Instruct* [7], *Qwen2.5-VL-72B-Instruct* [7], *QvQ-72B-Preview* [66], *InternVL2.5-38B* [15], *InternVL2.5-78B* [15], *InternVL3-38B* [94], *InternVL3-78B* [94], *LLaVA-v1.5-7B* [46], *LLaVA-OneVision-7B (SI)* [39], *ShareGPT4V* [9], *MiniCPM-o-2.6* [85], *GLM-4v-9B* [24], *Ovis2-8B* [53] and *mPLUG-Owl3-7B* [87], while close-source MLLMs consist of *GPT-4o* [32], *GPT-4o-mini*, *Kimi-latest* [65], *Doubaio-1.5-Vision-Pro-32k*, *Gemini-2.0-Pro* [64] and *Claude-3.7-Sonnet*. We further include two reinforcement-learning baselines built on *Qwen2.5-VL-7B-Instruct* [7] and *InternVL2.5-38B* [15], respectively, trained via our rule-based RL procedure on our supplementary training dataset. Fully supervised fine-tuning (SFT) experiments on the same datasets serve as controls to isolate the effect of RL optimization. All model hyperparameters, training regimes, and implementation details are provided in the Appendix.

LLM Evaluation Protocol. For language-only models, we generate an auxiliary image description using GPT-4o and prepend it to the question. Specifically, each question is formatted as “*Following is a detailed caption describing an image: [IMAGE DESCRIPTION]. Based on the provided description, select the best answer from the four options:*”. This combined prompt is fed directly into the target LLMs for inference.

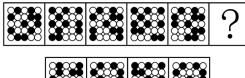
<p>Problem</p> <p>From the four options given, choose the most appropriate one to fill in the question mark so that it shows a certain regularity.</p>  <p>A B C D</p> <p>a</p>	<p>GPT-4o X</p> <p>To determine the correct answer, we need to identify the pattern in the sequence of images. Upon analyzing the sequence, each figure shows a progression where black and white dots alternate positions in a specific and regular way. Observations:</p> <ol style="list-style-type: none"> In each step, the black dots shift positions, creating a systematic pattern. The white and black dots alternate in a way that spreads uniformly across the grid. <p>By analyzing the progression in the sequence of grids, Option A is the correct choice as it aligns with the observed regularity.</p> <p>e</p>
<p>Reference Solution</p> <p>The elements are the same, so the position pattern is given priority. The figures in the question all have 25 squares, and the number of black blocks in the middle is the same. So the inner and outer circles are considered separately for translation. It is found that the position of the black block in the center of the inner circle remains unchanged, and the other two black blocks move 4 squares clockwise (counterclockwise) each time (it can also be seen as the positions of the two black blocks alternating on the left and right sides), and the outermost black block moves 1 square clockwise each time. Only option C is correct.</p> <p>b</p>	<p>QwenVL2.5 72B-Instruct X</p> <p>To determine the most appropriate option to fill in the question mark, we need to identify the pattern in the sequence of grids. Let's analyze the given grids step by step.</p> <p>The first grid has 10 black circles and 6 white circles.</p> <p>The second grid has 9 black circles and 7 white circles.</p> <p>The third grid has 8 black circles and 8 white circles.</p> <p>We can observe that the number of black circles decreases by 1 and the number of white circles increases by 1 in each subsequent grid. Therefore, the fourth grid should have 7 black circles and 9 white circles.</p> <p>Now, let's check the options:</p> <ul style="list-style-type: none"> A: This grid has 7 black circles and 9 white circles. B: This grid has 6 black circles and 10 white circles. C: This grid has 8 black circles and 8 white circles. D: This grid has 9 black circles and 7 white circles. <p>The grid that matches our pattern is option A, which has 7 black circles and 9 white circles.</p> <p>Thus, the most appropriate option to fill in the question mark is \boxed{A}.</p> <p>f</p>
<p>InternVL2.5 38B X</p> <p>According to the changing pattern of these graphs, we can observe that:</p> <p>The number of black circles in each graph is gradually decreasing, which allows us to analyze the options in the last row.</p> <p>Let's look at each option:</p> <ul style="list-style-type: none"> - The number of black circles in options A and B is more than in the previous graphs, which is obviously inconsistent with the pattern. - There are more black circles in option C than before, which is inconsistent with the previous pattern. - The number of black circles in option D is consistent with the previous change pattern. - Therefore, option D is the most suitable option. <p>c</p>	<p>InternVL2.5 38B + RL ✓</p> <p>To answer this question, we need to carefully observe the sequence of figures and find the pattern between them. First, let's look at the position and changes of the circles in each figure. The number of circles in each figure is the same, but their positions seem to be changing in some way. After careful analysis, we can see that the positions of the circles in the figures are moving in a certain order.</p> <ol style="list-style-type: none"> 1. In the first figure, the positions of the circles are relatively concentrated. 2. In the second figure, the positions of the circles change slightly and move slightly more than the first figure. 3. In the third figure, the positions of the circles change again and continue to move according to a certain pattern. 4. In the fourth figure, the positions of the circles are closer to the center. 5. In the fifth figure, the positions of the circles continue to move compared to the fourth figure. 6. In the sixth figure, the positions of the circles are different from the other figures, forming a new arrangement. <p>Through the analysis of the above figures, we can see that the current pattern of changes in the figures seems to be gradually forming a new arrangement. Therefore, the seventh figure should continue this change and form a new arrangement. According to this rule, the arrangement of option C is most consistent with the arrangement of the next figure. So, the answer is: C: C</p> <p>g</p>
<p>Image Description (GPT-4o)</p> <p>Structure of the Puzzle:</p> <ul style="list-style-type: none"> • The top row consists of five grids, each with a 6x6 arrangement of black and white circles. • The first four grids follow a specific pattern, while the fifth grid is missing (represented by a “?”). • The goal is to determine the pattern and select the correct option from four answer choices (A, B, C, D) shown below. <p>Observations About the Given Sequence:</p> <ol style="list-style-type: none"> 1. Each grid in the sequence contains a mix of black and white circles. 2. The arrangement of black circles seems to follow a transformation from one grid to the next. 3. By analyzing the first four grids, a rule or pattern governing the change can be identified. <p>Answer Choices (A, B, C, D):</p> <p>Each of the four answer choices consists of a 6x6 grid similar to the ones in the sequence. Only one of them correctly follows the pattern observed in the first four grids.</p> <p>d</p>	<p>03-mini X</p> <p>(Reasoned for 5 seconds.) To solve this puzzle, let's follow a step-by-step reasoning approach based on the observations about the sequence of grids.</p> <ol style="list-style-type: none"> 1. **First Step: Observe the Pattern in the First Four Grids** - We need to closely examine the black and white circles in each grid, paying attention to how the black circles evolve across each grid. - Look at the number of black circles, their positions, and any shifts or transformations from one grid to the next. 2. **Identify Possible Transformations** - The transformation could be based on the movement of black circles, changes in the arrangement (e.g., rotation, mirroring, shifting), or adding/removing black circles in a predictable manner. - Consider changes like: - Moving black circles from one position to another. - Adding or removing black circles. - Symmetry or patterns like shifting rows or columns. 3. **Predict the Missing Grid** - Using the transformations observed from the first four grids, hypothesize what the fifth grid should look like. - Think about the positions, number of black circles, and their movements. The fifth grid should be the natural progression based on the previous grids. 4. **Compare the Answer Choices** - Once you have a hypothesis for the missing grid, compare it against the four answer choices (A, B, C, D). - Choose the one that matches the predicted transformation. Based on this reasoning approach, let's assume that you've carefully analyzed the first four grids and identified a clear pattern in the movement and arrangement of the black circles. <p>Answer: \boxed{B}</p> <p>h</p>

Figure 6: **Solution examples generated by different models.** Reference solution and outputs generated by GPT-4o, Qwen2.5VLL-72B-Instruct, InternVL2.5-38B, and InternVL2.5-38B with RL. Additionally, the image caption and solution from LLMs (03-mini) are also illustrated.

Prompts Setting. We apply three distinct prompting paradigms to investigate model reasoning capabilities: 1) **Non-CoT prompt evaluation:** Models receive a concise instruction: “*Answer the question using a single word or phrase, following this format: Answer: \boxed{\$LETTER}\$*”. 2) **CoT prompt evaluation:** We prompt models to articulate intermediate reasoning steps: “*Solve the complex visual logical reasoning problem through step-by-step reasoning. Think about the reasoning process first and answer the question following this format: Answer: \boxed{\$LETTER\$}*”. 3) **Hint prompts evaluation:** Leveraging GPT-4o, we generate question-specific hints derived from the reference solutions. As shown in Figure 7, solution-related hints are provided alongside the CoT prompt to guide reasoning without revealing the final answer directly. Notably, unless otherwise specified, CoT prompt evaluation is employed by default for assessing model performance.

4.2 Overall Results

LLM Performance. Table 1 reports that all evaluated large language models attain rather low accuracy on VisuLogic. The best-performing LLM, *Qwen2.5-72B-Instruct*, reaches only 28.0%, while *GPT-4* and *Deepseek-R1* achieve 23.6% and 26.6%, respectively. These findings underscore that reasoning based solely on textual descriptions is insufficient to capture the rich visual information required by our benchmark and thus fails to resolve visual logical reasoning problems.

Table 1: **Cross-Modal performance with CoT prompts on VisuLogic.** The table shows the evaluation scores of baseline references, LLMs, and MLLMs, which illustrates a gap between humans’ and models’ capabilities. Top performers per category are **bolded**, with secondary leaders underlined.

Models	Overall	Quantity	Spatiality	Position	Attribute	Style	Other
References							
Human Random	51.4 24.9	45.3 25.7	52.7 25.4	71.1 22.7	50.0 23.4	47.5 24.3	44.2 26.1
Open Source LLM (MLLM Description→LLM)							
Deepseek-R1 [17] Qwen2.5-72B-Instruct [83] QwQ-32B [67]	26.6 28.0 22.8	<u>27.7</u> 30.2 24.6	<u>23.5</u> 24.4 20.1	24.0 27.5 <u>25.4</u>	27.8 <u>26.5</u> 19.0	<u>23.0</u> 26.8 20.7	35.0 <u>30.8</u> 24.0
Close Source LLM (MLLM Description→LLM)							
GPT-4 (20240613) [1] o3-mini (20250131) Gemini-2.0-Flash-Thinking (20250121) [64] Claude-3.7-Sonnet (20250219) Doubao-1.5-Pro-32k (20250115)	23.6 24.6 23.4 25.9 26.6	21.2 <u>27.8</u> <u>23.2</u> <u>25.9</u> 30.0	<u>22.5</u> 18.8 26.0 <u>22.5</u> <u>22.5</u>	21.3 24.5 16.9 25.0 <u>25.0</u>	25.6 21.7 17.1 28.0 <u>25.6</u>	23.3 <u>25.6</u> 21.1 33.3 30.6	35.2 28.4 <u>33.3</u> 30.6 30.0
Close Source MLLMs							
GPT-4o-mini (20240718) GPT-4o (20240806) [32] Kimi-latest [65] Doubao-1.5-Vision-Pro-32k (20250115) Gemini-2.0-Pro (20250205) [64] Claude-3.7-Sonnet (20250219)	24.3 26.3 25.9 28.1 28.0 24.8	27.2 <u>28.6</u> 24.9 28.1 29.7 22.7	23.4 24.7 29.4 23.8 <u>24.2</u> <u>27.3</u>	23.5 27.2 26.5 29.1 <u>27.9</u> <u>27.9</u>	18.3 26.8 <u>28.0</u> 25.1 30.5 <u>28.0</u>	<u>31.1</u> 20.0 16.7 32.1 22.2 22.2	16.7 25.9 26.9 35.0 <u>33.3</u> 22.2
Open Source MLLMs							
LLaVA-v1.5-7B [46] LLaVA-OneVision-7B (SI) [39] ShareGPT4V [9] MiniCPM-o-2.6 [85] GLM-4v-9B [24] Ovis2-8B [53] mPLUG-Owl3-7B-241101 [87] Qwen2.5-VL-7B-Instruct [7] Qwen2.5VL-72B-Instruct [7] QvQ-72B-Preview [66] InternVL2.5-38B [13] InternVL2.5-78B [13] InternVL3-38B [94] InternVL3-78B [94]	24.6 25.3 23.4 25.3 24.3 25.6 18.9 26.0 26.2 23.0 25.5 27.3 27.1 27.7	26.1 22.4 24.9 25.6 22.4 26.1 21.5 27.6 25.2 24.2 24.4 26.6 28.7 <u>27.7</u>	24.2 <u>27.3</u> 22.1 23.0 23.7 23.8 15.2 20.9 23.8 17.0 26.4 26.0 27.6 <u>27.6</u>	23.5 33.1 23.5 27.3 28.3 27.2 16.2 25.2 27.2 24.4 27.2 26.5 26.8 <u>26.8</u>	17.1 23.2 19.5 21.9 26.0 28.0 20.7 23.2 25.6 21.0 23.2 26.8 <u>31.1</u>	<u>31.1</u> 25.6 28.9 24.5 24.1 25.6 18.9 37.8 25.6 24.4 25.6 31.1 23.9 <u>32.3</u>	22.2 22.2 19.4 29.9 25.3 24.1 20.4 25.0 30.6 26.9 30.6 28.5 28.5 <u>32.3</u>
Qwen2.5-VL-7B-Instruct-SFT Qwen2.5-VL-7B-Instruct-RL InternVL2.5-38B-SFT InternVL2.5-38B-RL	25.5 28.0 27.9 31.1	24.4 26.6 30.6 31.2	26.4 33.8 29.4 <u>31.2</u>	<u>27.2</u> 29.4 20.6 26.5	23.2 23.2 25.6 26.5	25.6 18.9 30.0 30.5	26.9 <u>29.6</u> 25.0 38.9

MLLM Performance. As shown in Table 1, current multimodal LLMs also perform poorly on VisuLogic. The highest score is 28.1% by *Doubao-1.5-Vision-Pro-32k*, which remains a substantial 23.3 points below human performance. Advanced models such as *GPT-4o* and *Gemini-2.0-Pro* attain only 26.3% and 28.0%, respectively, revealing a marked gap between existing MLLMs and human-level visual reasoning. Overall, these results indicate that current MLLMs have serious deficiencies in visual reasoning and that significant advances are still required.

Effectiveness of CoT Prompts. Contrary to expectations, chain-of-thought (CoT) prompting yields minimal improvement in visual reasoning. As detailed in Table 2, *GPT-4o-mini* benefits most, with only a 1.2-point gain under CoT compared to direct-answer prompts; all other models exhibit gains below 1.0 point. We speculate that this limited effect likely stems from current CoT training being based only on pure-text corpora; future work should explore CoT techniques tailored to multimodal data to better support visual reasoning tasks.

Table 2: **Influence of Chain-of-Thought on model performance.** Positive value changes are highlighted in red, negative changes in green, and statistically insignificant variations (delta < 1%) are denoted in gray. MLLMs exhibit reasoning capability disparities, with some models demonstrating performance gains through chain-of-thought, while others reveal few effectiveness.

Models	CoT	Overall	Quantity	Spatiality	Position	Attribute	Style	Other
GPT-4o (20240806)	✓	26.3	28.6	24.7	27.2	26.8	20.0	25.9
	✗	26.0 _(-0.3)	26.9 _(-1.7)	24.2 _(-0.5)	26.5 _(-0.7)	23.2 _(-3.6)	24.0 _(+4.0)	29.6 _(+3.7)
Kimi-latest	✓	25.9	24.9	29.4	26.5	28.0	16.7	26.9
	✗	25.1 _(-0.8)	22.9 _(-2.0)	22.5 _(-6.9)	25.0 _(-1.5)	19.5 _(-7.5)	35.6 _(+18.9)	24.1 _(-2.8)
GPT-4o-mini (20240718)	✓	24.3	27.2	23.4	23.5	18.3	31.1	16.7
	✗	23.1 _(-1.2)	23.8 _(-3.4)	22.9 _(-0.5)	24.3 _(+0.8)	17.1 _(-1.2)	30.0 _(-1.1)	18.5 _(+1.8)
Qwen2.5-VL-Instruct-7B	✓	26.0	27.6	20.9	25.2	23.2	37.8	25.0
	✗	25.9 _(-0.1)	25.5 _(-2.1)	22.8 _(+1.9)	26.4 _(+1.2)	25.3 _(+2.1)	20.6 _(-17.2)	38.2 _(+13.2)
InternVL2.5-38B	✓	24.9	24.1	26.4	27.2	23.2	25.6	22.2
	✗	25.0 _(+0.1)	24.6 _(+0.5)	25.5 _(-0.9)	22.1 _(-5.1)	22.0 _(-1.2)	26.7 _(+1.1)	29.6 _(+7.4)

Table 3: **Influence of hint prompts on model performance.** MLLMs exhibit measurable performance enhancements with hint integration yet retain significant fidelity gaps against human performance. In comparison, humans achieve task mastery on VisuLogic, exposing critical gaps in models’ capacity. Value changes are color-coded with red indicating positive shifts and green denoting negative variations.

Models	Hint	Overall	Quantity	Spatiality	Position	Attribute	Style	Other
Human	✗	51.4	45.3	52.7	71.1	50.0	47.5	44.2
	✓	83.6 _(+32.2)	85.1 _(+39.8)	68.5 _(+15.8)	100.0 _(+28.9)	95.7 _(+45.7)	78.6 _(+31.1)	90.5 _(+46.3)
GPT-4o (20240806)	✗	26.3	28.6	24.7	27.2	26.8	20.0	25.9
	✓	30.0 _(+3.7)	25.4 _(-3.2)	31.5 _(+6.8)	29.2 _(+2.0)	28.6 _(+1.8)	30.8 _(+10.8)	42.9 _(+17.0)
Claude-3.7-Sonnet (20250219)	✗	24.8	22.7	27.3	27.9	28.0	22.2	22.2
	✓	33.5 _(+8.7)	37.3 _(+14.6)	33.3 _(+6.0)	37.5 _(+9.6)	23.8 _(-4.2)	15.4 _(-6.8)	38.1 _(+15.9)
Gemini-2.0-Pro (20250205)	✗	28.0	29.7	24.2	27.9	30.5	22.2	33.3
	✓	36.5 _(+8.5)	44.8 _(+15.1)	33.3 _(+9.1)	25.0 _(-2.9)	38.1 _(+7.6)	15.4 _(-6.8)	42.9 _(+9.6)
Doubaot-1.5-Vision-Pro-32k (20250115)	✗	28.1	28.1	23.8	29.1	25.1	32.1	35.0
	✓	37.0 _(+8.9)	46.3 _(+18.2)	25.9 _(+2.1)	54.2 _(+25.1)	33.3 _(+8.2)	23.1 _(-9.0)	28.6 _(-6.4)

Effectiveness of Hint Prompts. Table 3 shows that hint prompts can boost model performance—*Claude-3.7-Sonnet*, *Gemini-2.0-Pro*, and *Doubaot-1.5-Vision-Pro-32k* all improve by over 8 points, reaching accuracies above 35%. However, even with explicit guidance, models still fail to construct coherent, reliable reasoning chains. This suggests that simply augmenting training data with similar tasks is insufficient (which can help MLLMs come up with specific directions for solving the problem); future efforts must focus on enhancing the reliability and correctness of reasoning procedures of MLLMs to achieve more accurate reasoning inference.

Open-Source vs Close-Source. Table 1 further compares open- and closed-source models. The top open-source MLLM, *InternVL3-78B*, attains 27.7%, trailing the closed-source leader (*Doubaot-1.5-Vision-Pro-32k*, 28.1%) by only 0.4% points and outperforming other proprietary competitors such as *GPT-4o* and *Claude-3.7-Sonnet*. Overall, both open- and closed-source models exhibit uniformly low performance, highlighting a widespread neglect of visual reasoning objectives in current multimodal model training and data collection.

Behaviors of RL Trained models. As shown in Table 1, MLLMs with reinforcement learning optimization can yield significant improvements in visual reasoning performance. *Qwen2.5-VL-7B-Instruct-RL* attains 28.0%, a 2.0 percentage point boost over its non-RL counterpart. More strikingly, *InternVL2.5-38B-RL* reaches 31.1%, surpassing the original non-RL model by 5.6% and establishing a new state-of-the-art on VisuLogic. Furthermore, compared to supervised fine-tuning (SFT) on identical datasets, RL-enhanced models demonstrate substantially larger performance gains, underscoring the promise of targeted RL methods for advancing multimodal visual reasoning.

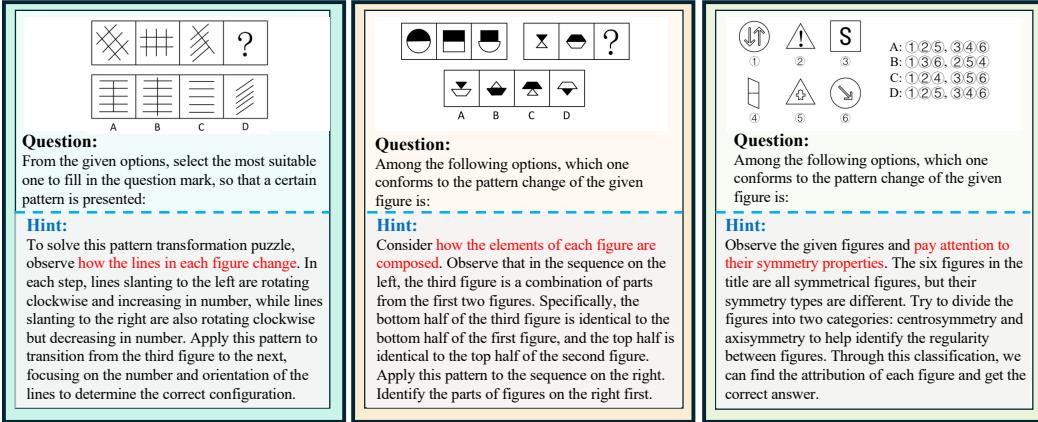


Figure 7: **Hint prompts visualization.** Hint prompts examples, which supply solution guidance for MLLMs, are shown in the image, with solution-critical elements highlighted in red.

4.3 Fine-grained Comparison

We systematically analyze model capabilities by examining error distributions across reasoning categories for different models. Figure 8 presents the error rates of LLMs, MLLMs, and human participants over six distinct reasoning categories.

Figure 8a reveals that LLMs struggle most with *Spatial Reasoning* questions, indicating that text-only descriptions are insufficient to infer three-dimensional structures or spatial transformations. In contrast, their performance on *Quantitative Reasoning* tasks is comparatively stronger, suggesting that quantitative relationships are more readily conveyed through language.

As shown in Figure 8b, *Stylistic Reasoning* presents the greatest difficulty for MLLMs, with error rates exceeding 75%—worse than random guessing (25% accuracy). This result underscores a fundamental limitation of current MLLM architectures in capturing subtle visual cues such as overlays, contours, and shape variations.

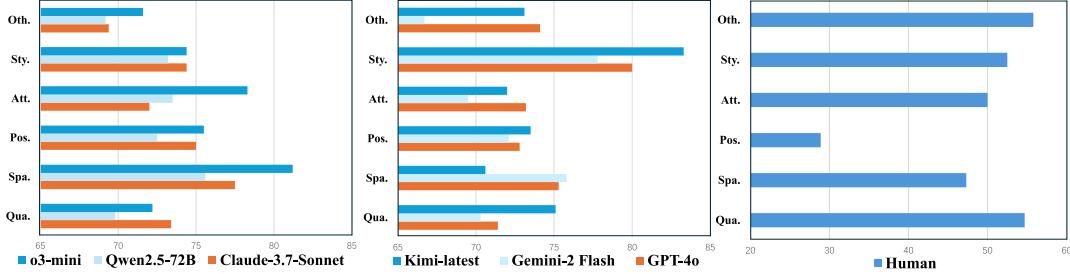
Figure 8c reveals that human error patterns form a distinct cluster, separate from both LLMs and MLLMs. Human participants maintain error rates below 30% on *Positional Reasoning* tasks, reflecting robust position-based visual inference. By contrast, both model classes struggle with positional reasoning, highlighting a fundamental divergence in visual–cognitive processes between humans and state-of-the-art large models.

4.4 Qualitative Analysis

LLM Failures. As shown in Figure 6(h), text-only LLMs that rely on externally generated image captions often omit critical visual details required for multi-step logical deduction—such as the counts, shapes, and progression patterns of the black and white dots in Figure 6(a). Consequently, their reasoning diverges from the correct solution and frequently yields hallucinations or irrelevant responses.

MLLM Failures. Figure 6 also presents cases in which MLLMs correctly describe static visual content yet fail to infer the evolving relationships among shapes, instead resorting to superficial cues like object counts. While these models can recognize individual shapes and tally items, they struggle to reason over inter-element relations, which limits their ability to solve complex visual-logic problems.

RL-Based Improvements. As illustrated in Figure 6(g), reinforcement learning (RL) fine-tuning encourages deeper, stepwise logical reasoning. The RL-enhanced model successfully captures state transitions (e.g., the movements of chess pieces in Figure 6(a)) and accurately predicts subsequent configurations. Moreover, it learns to iteratively revise intermediate hypotheses—akin to trial-and-error—until a coherent deduction emerges (see additional examples in the Appendix). These findings highlight the potential of RL methods to bolster performance on visual reasoning tasks.



(a) The error distribution of LLMs. (b) The error distribution of MLLMs. (c) The error distribution of humans.

Figure 8: **Error distribution analysis.** The figure demonstrates distinct error type allocations across Humans, LLMs and MLLMs, revealing differences among their cognition patterns.

5 Conclusion

In this paper, we present VisuLogic, a novel benchmark designed to evaluate the visual reasoning capabilities of Multi-modal Large Language Models (MLLMs). The benchmark consists of 1,000 vision-centric reasoning tasks distributed across six distinct categories. We conduct a comprehensive evaluation of several state-of-the-art LLMs and MLLMs on this benchmark and provide an in-depth analysis of their performance. Our findings reveal that even the most advanced models fall short of human performance, highlighting substantial opportunities for advancement in visual logical reasoning. Through the further experiments, we find that reinforcement learning (RL) is a promising approach for enhancing the vision reasoning capabilities of MLLMs. To promote further research and innovation, we have open-sourced the evaluation code, training scripts, and datasets associated with this work.

References

- [1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] W. U. Ahmad, S. Narendhiran, S. Majumdar, A. Ficek, S. Jain, J. Huang, V. Noroozi, and B. Ginsburg. Opencodereasoning: Advancing data distillation for competitive coding. *arXiv preprint arXiv:2504.01943*, 2025.
- [3] S. N. Akter, S. Lee, Y. Chang, Y. Bisk, and E. Nyberg. Visreas: Complex visual reasoning with unanswerable questions, 2024.
- [4] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.
- [5] H. Bai, Y. Zhou, L. E. Li, S. Levine, and A. Kumar. Digi-q: Learning q-value functions for training device-control agents. *arXiv preprint arXiv:2502.15760*, 2025.
- [6] J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, and J. Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- [7] S. Bai, K. Chen, X. Liu, J. Wang, W. Ge, S. Song, K. Dang, P. Wang, S. Wang, J. Tang, H. Zhong, Y. Zhu, M. Yang, Z. Li, J. Wan, P. Wang, W. Ding, Z. Fu, Y. Xu, J. Ye, X. Zhang, T. Xie, Z. Cheng, H. Zhang, Z. Yang, H. Xu, and J. Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- [8] G. Chen, M. Wang, Y. Yang, K. Yu, L. Yuan, and Y. Yue. Pointgpt: Auto-regressively generative pre-training from point clouds. *Advances in Neural Information Processing Systems*, 36:29667–29679, 2023.
- [9] L. Chen, J. Li, X. Dong, P. Zhang, C. He, J. Wang, F. Zhao, and D. Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023.
- [10] L. Chen, L. Li, H. Zhao, Y. Song, and Vinci. R1-v: Reinforcing super generalization ability in vision-language models with less than \$3. <https://github.com/Deep-Agent/R1-V>, 2025. Accessed: 2025-02-02.
- [11] S. Chen, H. Li, Q. Wang, Z. Zhao, M. Sun, X. Zhu, and J. Liu. Vast: A vision-audio-subtitle-text omni-modality foundation model and dataset. *Advances in Neural Information Processing Systems*, 36:72842–72866, 2023.
- [12] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollar, and C. L. Zitnick. Microsoft coco captions: Data collection and evaluation server, 2015.
- [13] Z. Chen, W. Wang, Y. Cao, Y. Liu, Z. Gao, E. Cui, J. Zhu, S. Ye, H. Tian, Z. Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.
- [14] Z. Chen, W. Wang, H. Tian, S. Ye, Z. Gao, E. Cui, W. Tong, K. Hu, J. Luo, Z. Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024.
- [15] Z. Chen, J. Wu, W. Wang, W. Su, G. Chen, S. Xing, M. Zhong, Q. Zhang, X. Zhu, L. Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24185–24198, 2024.
- [16] S. Cheng, Z. Guo, J. Wu, K. Fang, P. Li, H. Liu, and Y. Liu. Egothink: Evaluating first-person perspective thinking capability of vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14291–14302, 2024.

- [17] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.
- [18] A. Défossez, L. Mazaré, M. Orsini, A. Royer, P. Pérez, H. Jégou, E. Grave, and N. Zeghidour. Moshi: a speech-text foundation model for real-time dialogue. *arXiv preprint arXiv:2410.00037*, 2024.
- [19] H. Dong, J. Li, B. Wu, J. Wang, Y. Zhang, and H. Guo. Benchmarking and improving detail image caption. *arXiv preprint arXiv:2405.19092*, 2024.
- [20] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [21] Q. Fang, S. Guo, Y. Zhou, Z. Ma, S. Zhang, and Y. Feng. Llama-omni: Seamless speech interaction with large language models. *arXiv preprint arXiv:2409.06666*, 2024.
- [22] J. Feng, R. Xu, J. Hao, H. Sharma, Y. Shen, D. Zhao, and W. Chen. Language models can be logical solvers. *arXiv preprint arXiv:2311.06158*, 2023.
- [23] Z. Gao, Z. Chen, E. Cui, Y. Ren, W. Wang, J. Zhu, H. Tian, S. Ye, J. He, X. Zhu, et al. Mini-internvl: a flexible-transfer pocket multi-modal model with 5% parameters and 90% performance. *Visual Intelligence*, 2(1):1–17, 2024.
- [24] T. GLM, A. Zeng, B. Xu, B. Wang, C. Zhang, D. Yin, D. Rojas, G. Feng, H. Zhao, H. Lai, H. Yu, H. Wang, J. Sun, J. Zhang, J. Cheng, J. Gui, J. Tang, J. Zhang, J. Li, L. Zhao, L. Wu, L. Zhong, M. Liu, M. Huang, P. Zhang, Q. Zheng, R. Lu, S. Duan, S. Zhang, S. Cao, S. Yang, W. L. Tam, W. Zhao, X. Liu, X. Xia, X. Zhang, X. Gu, X. Lv, X. Liu, X. Liu, X. Yang, X. Song, X. Zhang, Y. An, Y. Xu, Y. Niu, Y. Yang, Y. Li, Y. Bai, Y. Dong, Z. Qi, Z. Wang, Z. Yang, Z. Du, Z. Hou, and Z. Wang. Chatglm: A family of large language models from glm-130b to glm-4 all tools, 2024.
- [25] B. Goertzel and C. Pennachin. *Artificial general intelligence*, volume 2. Springer, 2007.
- [26] Z. Guo, R. Zhang, X. Zhu, Y. Tang, X. Ma, J. Han, K. Chen, P. Gao, X. Li, H. Li, et al. Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. *arXiv preprint arXiv:2309.00615*, 2023.
- [27] H. Gupta, S. Verma, U. Anantheswaran, K. Scaria, M. Parmar, S. Mishra, and C. Baral. Polymath: A challenging multi-modal mathematical reasoning benchmark, 2024.
- [28] X. Han, Q. You, Y. Liu, W. Chen, H. Zheng, K. Mrini, X. Lin, Y. Wang, B. Zhai, J. Yuan, et al. Infimm-eval: Complex open-ended reasoning evaluation for multi-modal large language models. *arXiv preprint arXiv:2311.11567*, 2023.
- [29] Y. Hao, J. Gu, H. W. Wang, L. Li, Z. Yang, L. Wang, and Y. Cheng. Can mllms reason in multimodality? emma: An enhanced multimodal reasoning benchmark. *arXiv preprint arXiv:2501.05444*, 2025.
- [30] C. He, R. Luo, Y. Bai, S. Hu, Z. L. Thai, J. Shen, J. Hu, X. Han, Y. Huang, Y. Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint arXiv:2402.14008*, 2024.
- [31] D. Huang, Q. Bu, Y. Qing, and H. Cui. Codecto: Tackling code syntax errors in cot reasoning for code generation. *arXiv preprint arXiv:2308.08784*, 2023.
- [32] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- [33] A. Jaech, A. Kalai, A. Lerer, A. Richardson, A. El-Kishky, A. Low, A. Helyar, A. Madry, A. Beutel, A. Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- [34] X. Jiang, Y. Dong, L. Wang, Z. Fang, Q. Shang, G. Li, Z. Jin, and W. Jiao. Self-planning code generation with large language models. *ACM Transactions on Software Engineering and Methodology*, 33(7):1–30, 2024.

- [35] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, M. Suleyman, and A. Zisserman. The kinetics human action video dataset, 2017.
- [36] L. Ke, W. Pei, R. Li, X. Shen, and Y.-W. Tai. Reflective decoding network for image captioning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8888–8897, 2019.
- [37] W. Kool, H. van Hoof, and M. Welling. Buy 4 reinforce samples, get a baseline for free! 2019.
- [38] X. Lai, Z. Tian, Y. Chen, Y. Li, Y. Yuan, S. Liu, and J. Jia. Lisa: Reasoning segmentation via large language model, 2024.
- [39] B. Li, Y. Zhang, D. Guo, R. Zhang, F. Li, H. Zhang, K. Zhang, P. Zhang, Y. Li, Z. Liu, and C. Li. Llava-onevision: Easy visual task transfer, 2024.
- [40] J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023.
- [41] J. Li, D. Li, C. Xiong, and S. Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pages 12888–12900. PMLR, 2022.
- [42] J. Li, G. Li, Y. Li, and Z. Jin. Structured chain-of-thought prompting for code generation. *ACM Transactions on Software Engineering and Methodology*, 34(2):1–23, 2025.
- [43] J. Li, W. Lu, H. Fei, M. Luo, M. Dai, M. Xia, Y. Jin, Z. Gan, D. Qi, C. Fu, Y. Tai, W. Yang, Y. Wang, and C. Wang. A survey on benchmarks of multimodal large language models, 2024.
- [44] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollár. Microsoft coco: Common objects in context, 2015.
- [45] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.
- [46] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning, 2023.
- [47] H. Liu, Z. Teng, L. Cui, C. Zhang, Q. Zhou, and Y. Zhang. Logicot: Logical chain-of-thought instruction-tuning. *arXiv preprint arXiv:2305.12147*, 2023.
- [48] X. Liu, H. Yu, H. Zhang, Y. Xu, X. Lei, H. Lai, Y. Gu, H. Ding, K. Men, K. Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*, 2023.
- [49] Y. Liu, Z. Li, M. Huang, B. Yang, W. Yu, C. Li, X.-C. Yin, C.-L. Liu, L. Jin, and X. Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. *Science China Information Sciences*, 67(12):220102, 2024.
- [50] Z. Liu, Z. Sun, Y. Zang, X. Dong, Y. Cao, H. Duan, D. Lin, and J. Wang. Visual-rft: Visual reinforcement fine-tuning. *arXiv preprint arXiv:2503.01785*, 2025.
- [51] Z. Liu, Y. Zhang, F. Liu, C. Zhang, Y. Sun, and J. Wang. Othink-mrl: Stimulating multimodal generalized reasoning capabilities through dynamic reinforcement learning. *arXiv preprint arXiv:2503.16081*, 2025.
- [52] P. Lu, H. Bansal, T. Xia, J. Liu, C. Li, H. Hajishirzi, H. Cheng, K.-W. Chang, M. Galley, and J. Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- [53] S. Lu, Y. Li, Q.-G. Chen, Z. Xu, W. Luo, K. Zhang, and H.-J. Ye. Ovis: Structural embedding alignment for multimodal large language model. *arXiv:2405.20797*, 2024.
- [54] K. Mangalam, R. Akshulakov, and J. Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *Advances in Neural Information Processing Systems*, 36:46212–46244, 2023.

- [55] J. Mao, J. Huang, A. Toshev, O. Camburu, A. L. Yuille, and K. Murphy. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 11–20, 2016.
- [56] A. Masry, D. X. Long, J. Q. Tan, S. Joty, and E. Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022.
- [57] M. Mathew, D. Karatzas, and C. Jawahar. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209, 2021.
- [58] F. Meng, L. Du, Z. Liu, Z. Zhou, Q. Lu, D. Fu, B. Shi, W. Wang, J. He, K. Zhang, et al. Mm-eureka: Exploring visual aha moment with rule-based large-scale reinforcement learning. *arXiv preprint arXiv:2503.07365*, 2025.
- [59] T. Nguyen, S. Y. Gadre, G. Ilharco, S. Oh, and L. Schmidt. Improving multimodal datasets with image captioning. *Advances in Neural Information Processing Systems*, 36:22047–22069, 2023.
- [60] R. Niu, J. Li, S. Wang, Y. Fu, X. Hu, X. Leng, H. Kong, Y. Chang, and Q. Wang. Screenagent: A vision language model-driven computer control agent. *arXiv preprint arXiv:2402.07945*, 2024.
- [61] Y. Peng, G. Zhang, M. Zhang, Z. You, J. Liu, Q. Zhu, K. Yang, X. Xu, X. Geng, and X. Yang. Lmm-r1: Empowering 3b lmms with strong reasoning abilities through two-stage rule-based rl. *arXiv preprint arXiv:2503.07536*, 2025.
- [62] R. Qiao, Q. Tan, G. Dong, M. Wu, C. Sun, X. Song, Z. GongQue, S. Lei, Z. Wei, M. Zhang, et al. We-math: Does your large multimodal model achieve human-like mathematical reasoning? *arXiv preprint arXiv:2407.01284*, 2024.
- [63] H. Shen, Z. Zhang, K. Zhao, Q. Zhang, R. Xu, and T. Zhao. Vlm-r1: A stable and generalizable r1-style large vision-language model. <https://github.com/om-ai-lab/VLM-R1>, 2025. Accessed: 2025-02-15.
- [64] G. Team, R. Anil, S. Borgeaud, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, K. Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [65] K. Team, A. Du, B. Gao, B. Xing, C. Jiang, C. Chen, C. Li, C. Xiao, C. Du, C. Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.
- [66] Q. Team. Qvq: To see the world with wisdom, December 2024.
- [67] Q. Team. Qwq-32b: Embracing the power of reinforcement learning, March 2025.
- [68] Y. Wan, W. Wang, Y. Yang, Y. Yuan, J.-t. Huang, P. He, W. Jiao, and M. R. Lyu. Logicasker: Evaluating and improving the logical reasoning ability of large language models. *arXiv preprint arXiv:2401.00757*, 2024.
- [69] K. Wang, J. Pan, W. Shi, Z. Lu, H. Ren, A. Zhou, M. Zhan, and H. Li. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in Neural Information Processing Systems*, 37:95095–95169, 2024.
- [70] P. Wang, S. Bai, S. Tan, S. Wang, Z. Fan, J. Bai, K. Chen, X. Liu, J. Wang, W. Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- [71] S. Wang, D. Kim, A. Taalimi, C. Sun, and W. Kuo. Learning visual grounding from generative vision and language model. In *2025 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 8057–8067. IEEE, 2025.

- [72] W. Wang, Z. Gao, L. Chen, Z. Chen, J. Zhu, X. Zhao, Y. Liu, Y. Cao, S. Ye, X. Zhu, et al. Visualprm: An effective process reward model for multimodal reasoning. *arXiv preprint arXiv:2503.10291*, 2025.
- [73] X. Wang and P. Peng. Open-r1-video. <https://github.com/Wang-Xiaodong1899/Open-R1-Video>, 2025.
- [74] Y. Wang, W. Chen, X. Han, X. Lin, H. Zhao, Y. Liu, B. Zhai, J. Yuan, Q. You, and H. Yang. Exploring the reasoning abilities of multimodal large language models (mllms): A comprehensive survey on emerging trends in multimodal reasoning, 2024.
- [75] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [76] Y. Xiao, E. Sun, T. Liu, and W. Wang. Logicvista: Multimodal llm logical reasoning benchmark in visual contexts. *arXiv preprint arXiv:2407.04973*, 2024.
- [77] Z. Xie and C. Wu. Mini-omni: Language models can hear, talk while thinking in streaming, 2024. URL <https://arxiv.org/abs/2408.16725>.
- [78] D. Xu, Z. Zhao, J. Xiao, F. Wu, H. Zhang, X. He, and Y. Zhuang. Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 1645–1653, 2017.
- [79] F. Xu, Z. Wu, Q. Sun, S. Ren, F. Yuan, S. Yuan, Q. Lin, Y. Qiao, and J. Liu. Symbol-llm: Towards foundational symbol-centric interface for large language models. *arXiv preprint arXiv:2311.09278*, 2023.
- [80] R. Xu, Z. Huang, T. Wang, Y. Chen, J. Pang, and D. Lin. Vlm-grounder: A vlm agent for zero-shot 3d visual grounding. *arXiv preprint arXiv:2410.13860*, 2024.
- [81] Y. Xu, X. Liu, X. Liu, Z. Hou, Y. Li, X. Zhang, Z. Wang, A. Zeng, Z. Du, W. Zhao, et al. Chatglm-math: Improving math problem-solving in large language models with a self-critique pipeline. *arXiv preprint arXiv:2404.02893*, 2024.
- [82] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, et al. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [83] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, H. Lin, J. Yang, J. Tu, J. Zhang, J. Yang, J. Yang, J. Zhou, J. Lin, K. Dang, K. Lu, K. Bao, K. Yang, L. Yu, M. Li, M. Xue, P. Zhang, Q. Zhu, R. Men, R. Lin, T. Li, T. Xia, X. Ren, X. Ren, Y. Fan, Y. Su, Y. Zhang, Y. Wan, Y. Liu, Z. Cui, Z. Zhang, and Z. Qiu. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [84] Y. Yang, X. He, H. Pan, X. Jiang, Y. Deng, X. Yang, H. Lu, D. Yin, F. Rao, M. Zhu, B. Zhang, and W. Chen. R1-onevision: Advancing generalized multimodal reasoning through cross-modal formalization. *arXiv preprint arXiv:2503.10615*, 2025.
- [85] Y. Yao, T. Yu, A. Zhang, C. Wang, J. Cui, H. Zhu, T. Cai, H. Li, W. Zhao, Z. He, et al. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv preprint arXiv:2408.01800*, 2024.
- [86] J. Ye, G. Li, S. Gao, C. Huang, Y. Wu, S. Li, X. Fan, S. Dou, Q. Zhang, T. Gui, et al. Tooleyes: fine-grained evaluation for tool learning capabilities of large language models in real-world scenarios. *arXiv preprint arXiv:2401.00741*, 2024.
- [87] J. Ye, H. Xu, H. Liu, A. Hu, M. Yan, Q. Qian, J. Zhang, F. Huang, and J. Zhou. mplug-owl3: Towards long image-sequence understanding in multi-modal large language models, 2024.
- [88] L. Yu, P. Poirson, S. Yang, A. C. Berg, and T. L. Berg. Modeling context in referring expressions. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part II* 14, pages 69–85. Springer, 2016.

- [89] X. Yue, Y. Ni, K. Zhang, T. Zheng, R. Liu, G. Zhang, S. Stevens, D. Jiang, W. Ren, Y. Sun, C. Wei, B. Yu, R. Yuan, R. Sun, M. Yin, B. Zheng, Z. Yang, Y. Liu, W. Huang, H. Sun, Y. Su, and W. Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of CVPR*, 2024.
- [90] C. Zhang, F. Gao, B. Jia, Y. Zhu, and S.-C. Zhu. Raven: A dataset for relational and analogical visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [91] R. Zhang, D. Jiang, Y. Zhang, H. Lin, Z. Guo, P. Qiu, A. Zhou, P. Lu, K.-W. Chang, Y. Qiao, et al. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? In *European Conference on Computer Vision*, pages 169–186. Springer, 2024.
- [92] Q. Zhao, S. Wang, C. Zhang, C. Fu, M. Q. Do, N. Agarwal, K. Lee, and C. Sun. Antgpt: Can large language models help long-term action anticipation from videos? *arXiv preprint arXiv:2307.16368*, 2023.
- [93] D. Zhu, J. Chen, X. Shen, X. Li, and M. Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- [94] J. Zhu, W. Wang, Z. Chen, Z. Liu, S. Ye, L. Gu, Y. Duan, H. Tian, W. Su, J. Shao, Z. Gao, E. Cui, Y. Cao, Y. Liu, X. Wei, H. Zhang, H. Wang, W. Xu, H. Li, J. Wang, D. Chen, S. Li, Y. He, T. Jiang, J. Luo, Y. Wang, C. He, B. Shi, X. Zhang, W. Shao, J. He, Y. Xiong, W. Qu, P. Sun, P. Jiao, H. Lv, L. Wu, K. Zhang, H. Deng, J. Ge, K. Chen, L. Wang, M. Dou, L. Lu, X. Zhu, T. Lu, D. Lin, Y. Qiao, J. Dai, and W. Wang. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models, 2025.

A Overview of the Appendix

In the appendix, we provide additional information regarding the benchmark section and the experiment section. In Section B, we analyze the distribution characteristics of the data through statistical analysis and include examples of questions from different categories for reference. In the Section C, we first elaborate on the experimental settings for the large language model (LLM) experiments and provide details on the corresponding captions. We also present the outputs of different models on various types of questions for comparison. Then we describe our data-generation details and prompt examples for hint prompts experiment. Additionally, we supplement the details of our supervised fine-tuning (SFT) and reinforcement learning (RL) experiments, including examples for better illustration.

B Benchmark Analysis

B.1 Statistical analysis

As shown in Figure 10, the text length of our benchmark is mostly concentrated around 40 tokens (we calculated these statistics using the tokenizers of two open-source models). We also analyze the distribution of image sizes, as shown in Figure 9. The image widths range from 200 to 700 pixels, with an average of 592.3 pixels, while the heights range from 90 to 825 pixels, with an average of 327.9 pixels.

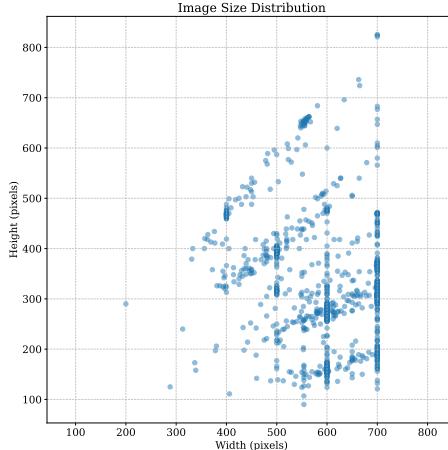


Figure 9: **Image Size Distribution.** The size of the image is basically limited to within the same order of magnitude.

B.2 More Questions of VisuLogic

To provide a thoroughly presentation of our benchmark, we include more questions from different categories in the Figure 11 and Figure 12. For the majority of questions in the benchmark, both the information and answer options are fully contained within the images. Only a very small fraction (less than 1%) of the questions, as illustrated by the question in the top-left corner of the ‘Other’ section in Figure 12, require additional textual context beyond what is shown in the images.

C Evaluation

C.1 Evaluation of LLMs

Caption generation for LLMs Evaluation. In our experiment, we employ large language models (LLMs) for comparative analysis. Specifically, when setting up the LLM-based experiment, we

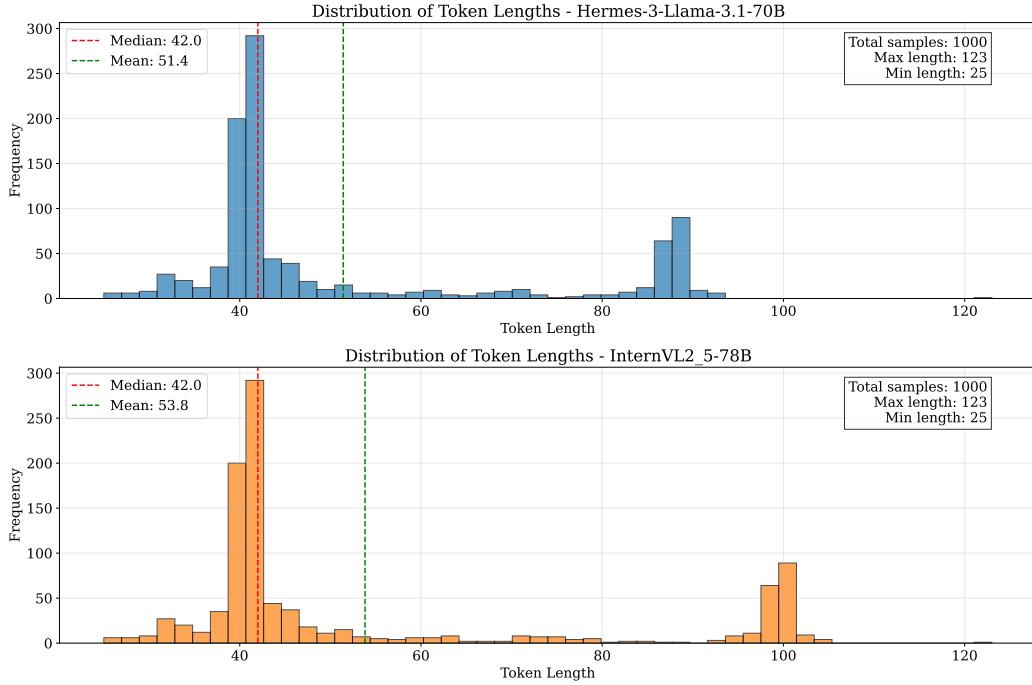


Figure 10: Distribution of text token length in VisuLogic

initially utilize GPT-4o to generate captions for images with the following prompt: There is a question about the image or figure. The prompt is as follow: *Please describe the fine-grained content of the image or figure based on this question, including scenes, objects, relationships, and any text present. Please note that you do not need to answer this question directly, just describe the information of this picture.* Additional examples of generated image captions are presented in Figure 14 and Figure 15.

More Examples of Captions. We provide additional image captions for six categories, as illustrated in Figures 14 and 15. Even SOTA MLLM (GPT-4o) encounters difficulties in accurately describing the details of images from VisuLogic.

C.2 More Solutions from Models

We provide more solutions generated from different LLMs/MLLMs on our benchmark, as shown in Figure 16, Figure 17 and Figure 18. For the majority of questions, almost all models fail to provide solutions meeting expectation, even when the final answer is correct.

C.3 Hint Prompts Evaluation Details

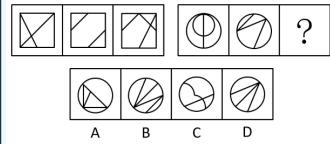
We first generate hint prompts with GPT-4o, combining reference solutions with question data as inputs. All outputs undergo manual validation to prevent solution leakage. More examples are shown in 19. After that, we input the hint prompts along with the same CoT prompt as CoT experiments (“*Solve the complex visual logical reasoning problem through step-by-step reasoning. Think about the reasoning process first and answer the question following this format: Answer: \boxed{\\$LETTER}.*”) to MLLMs.

C.4 RL Experiments

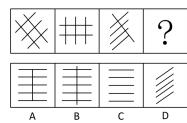
Comparative SFT Experiment. To verify the effectiveness of RL method, we arrange the comparative SFT experiments on the same dataset as RL experiments. The instruction consists of questions and CoT prompts, and the responses are formatted direct answers.

Quantitative Reasoning

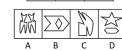
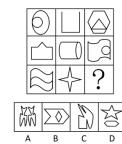
From the four given options, select the most suitable one to fill in the question mark:



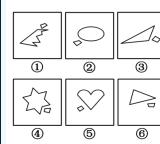
From the given four options, select the most suitable one to fill in the question mark.



From the given four options, choose the most suitable one to fill in the question mark:



Divide the following figures into two categories, ensuring that each category shares common patterns.



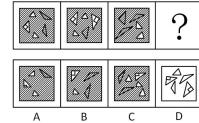
A: ①②⑥, ③④⑤

B: ①④⑤, ②③⑥

C: ①③⑥, ②④⑤

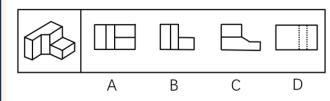
D: ①③⑤, ②④⑥

From the given four options, select the most appropriate one to fill in the question mark.

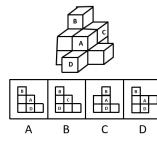


Spatial Reasoning

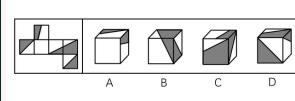
From any angle, which one on the right is not a view of the three-dimensional shape given on the left?



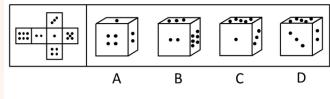
Which of the following four options is the front elevation view of the figure in the question:



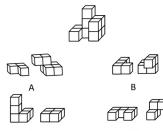
The image on the left shows the net of a cube. Which of the options on the right can be formed by folding it?



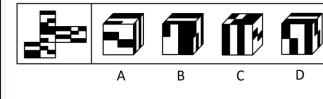
The diagram on the left shows the net of a cube's outer surface. Please identify the incorrect option.



Choose the option from the four given choices that cannot assemble into the 3D shape shown in the question.

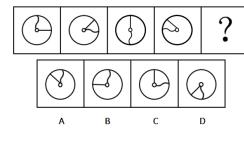


The left shows the net of a cube's outer surface. Which option on the right can be folded into it?

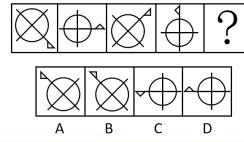


Positional Reasoning

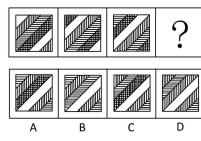
From the four given options, select the most suitable one to fill in the question mark:



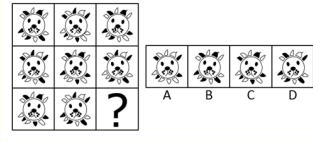
The option that best matches the given pattern is ().



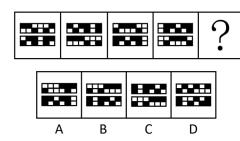
From the four given options, choose the most suitable one to fill in the question mark to create a certain pattern:



From the given four options, choose the most suitable one to fill in the question mark



From the four given options, choose the most appropriate one to fill in the question mark



From the given four options, choose the most suitable one to present a certain regularity:

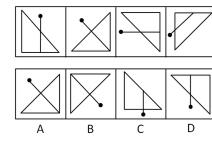
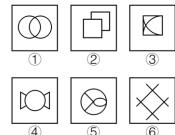


Figure 11: More examples in VisuLogic of Quantitative Reasoning, Spatial Reasoning, Positional Reasoning.

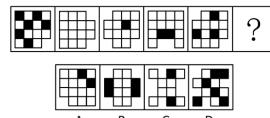
Attribute Reasoning

Divide the six figures below into two categories, so that each category of figures has its own common patterns.



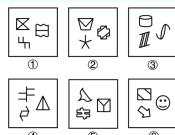
- A: ①②③, ④⑤⑥
B: ①②⑤, ③④⑥
C: ①③⑤, ②④⑥
D: ①④⑥, ②③⑤

From the four given options, choose the most suitable one to fill in the question mark.



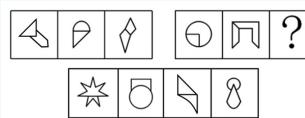
A B C D

Divide the six figures below into two categories, so that each category of figures has its own common patterns.



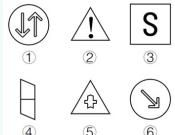
- A: ①②③, ④⑤⑥
B: ①②⑥, ③④⑤
C: ①③④, ②⑤⑥
D: ①④⑤, ②③⑥

From the given four options, select the most suitable one to fill in the question mark.



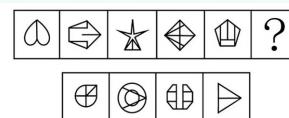
A B C D

Divide the following figures into two categories, ensuring that each category shares common patterns.



- A: ①②⑤, ③④⑥
B: ①③⑥, ②④⑤
C: ①②④, ③⑤⑥
D: ①③④, ②⑤⑥

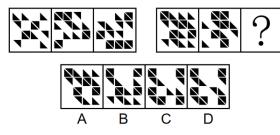
From the given four options, select the most appropriate one to fill in the question mark.



A B C D

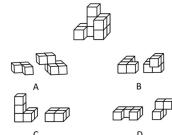
Stylistic Reasoning

From the four given options, choose the most suitable one to fill in the question mark.



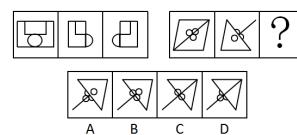
A B C D

Among the following options, the figure formed by overlapping the given figures is:



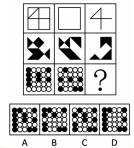
A B C D

Choose the most suitable option from the four given choices to fill in the question mark.



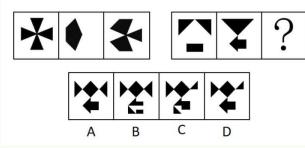
A B C D

From the four given options, choose the most suitable one to fill in the question mark.



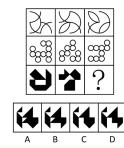
A B C D

From the four given options, choose the most suitable one to fill in the question mark.



A B C D

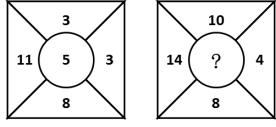
From the four given options, choose the most appropriate one to fill in the question mark.



A B C D

Other

From the four given options, choose the most suitable one to fill in the question mark: A: 2 B: 6 C: 4 D: 3

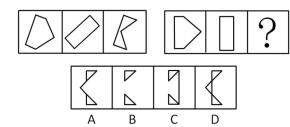


Choose the most suitable option from the four given choices to fill in the question mark.



A B C D

From the given four options, choose the most appropriate one to fill in the question mark.



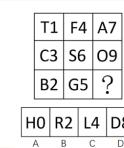
A B C D

Select the most suitable option from the four given choices to fill in the question mark.



A B C D

From the four given options, choose the most suitable one to fill in the question mark.



A B C D

Figure 12: More examples in VisuLogic of Attribute Reasoning, Stylistic Reasoning, and Other.

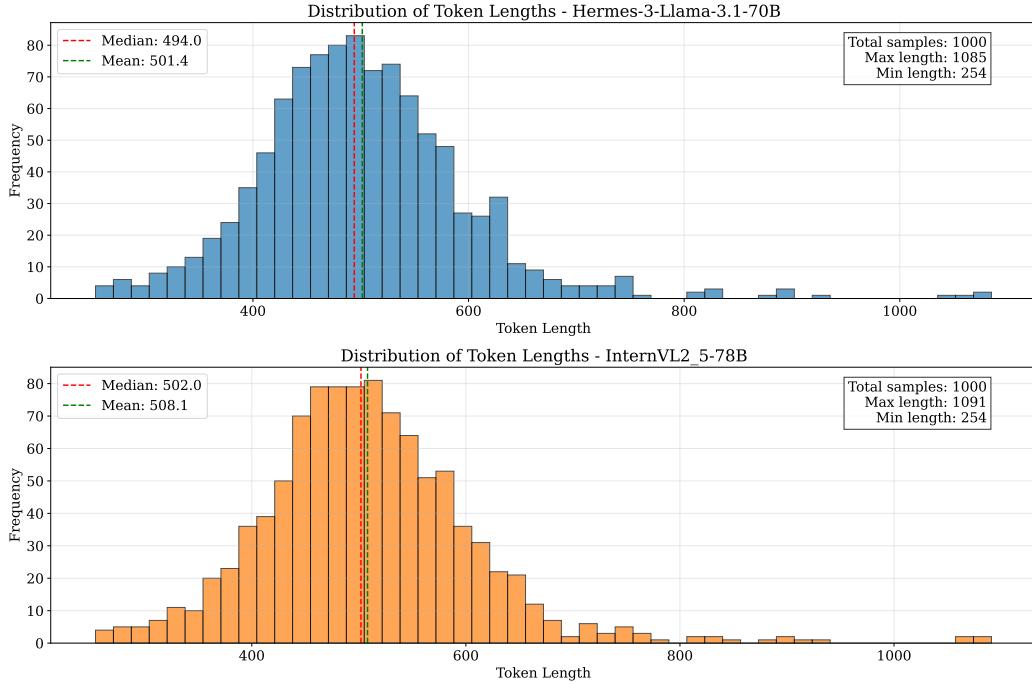


Figure 13: Distribution of benchmark with image caption.

RL Algorithm. We employ REINFORCE Leave-One-Out (RLOO) [37] in our reinforcement learning training phase. As a critic-model-free algorithm, rloo is at a low computational cost while maintaining more robustness to noise and KL constraints.

Reward Modeling. Inspired by Deepseek-R1 [17], we design our rule-based reward system that mainly consists of two types of rewards:

1. **Format rewards:** To clarify model's outputs, we design a format rule that forces model to put its thinking process between '<think>' and '</think>' tags and put its final answer between '<answer>' and '</answer>' tags. Regular expression is applied to judge whether outputs conform to the format rule.
2. **Accuracy rewards:** The accuracy reward is decided by the response's correctness. The model should generate the response in right format, then the answer is extracted and judged whether it is matched to the correct option.

Hyperparameter settings. Our two RL models are trained with the hyperparameter configuration detailed in Table 5. And the hyperparameters used in sft training stage are listed in Table 4.

Other Details. The training environment consists of CentOS Linux release 7.6.1810 operating system with CUDA 12.1. For Qwen2.5-VL-7B-Instruct-RL, we train for 40 steps on 1*8 A800 GPUS and for InternVL2.5-38B-RL we train for 60 steps on 6*8 A800 GPUS.

C.5 RL Experiments Examples

Figures 20, 21, 22, 23, 24 and 25 demonstrate qualitative differences in model outputs between baseline and RL optimized models. It illustrates reinforcement learning (RL) training enables the model to perform fundamental-level analysis of reasoning tasks embedded in graphical representations.

IMAGE DESCRIPTION:The image presents a visual pattern recognition puzzle consisting of a sequence of geometric shapes organized in two rows.

The objective is to determine the missing shape in the sequence from the given options.

Top Row (Sequence of Shapes)

- First Shape**:** - Two identical triangles, one upright and one upside down, overlap each other perfectly to form a star-like shape with a total of six outward-pointing triangular tips.
- The outlines of the triangles intersect, forming an internal hexagon.

Second Shape:** - Two overlapping triangles, one upright and one upside down, slightly offset and not perfectly aligned, creating a symmetric "X" formation at the center.

- The overlap forms a six-pointed star with more distinct and enhanced intersection lines visible inside.

Third Shape:** - Two overlapping triangles, one upright and one upside down, arranged symmetrically to create a six-pointed star.

- The internal lines are more connected, forming a central hexagon similar in style to the first shape, with additional internal lines from the triangle overlaps.

Fourth Shape:** - Two overlapping parallelograms, forming a shape that looks like a hexagon outlined in a bold, overlapping manner.

- The parallelograms intersect to create a pattern with sharp angles and straight lines defining internal quadrilateral shapes.

Fifth Shape:** - A centralized octagon with an inner pattern of lines creating a star shape, with the corners of the star connecting each side of the octagon from the center.

Missing Shape (Question Mark):** - This is the space where the next shape in the sequence would logically fit, represented by a question mark.

Bottom Row (Answer Options)

- Option A**:** - Two vertical rectangles with a central pointed triangle merging at the top to form an arrowhead shape.
- The base of the rectangles is flat, with an internal vertical line separating them.

Option B:** - Two triangles placed side by side with a small gap between them.

- The triangles are upright and identical in size and shape.

Option C:** - Two parallelograms overlapping each other slightly, tilted sideways.

- The overlapping creates an interesting central intersection with slanting lines.

Option D:** - An octagon with a central star pattern formed by crossing diagonal lines.

- The star is symmetrical and fits perfectly within the octagon's outline.

Task

The objective is to determine which of the four options (A, B, C, or D) logically completes the sequence in the top row.] From the four given options, choose the most suitable one to fill in the question mark, so that it presents a certain regularity:

A: A\nB: B\nC: C\nD: D

IMAGE DESCRIPTION:The image is a 3x3 grid puzzle with the task of identifying the missing symbol in the bottom right cell. Below the grid are four answer options labeled A, B, C, and D.

Grid Description

First Row:

- Cell 1**:** A circle with a diagonal line running from the top left to the bottom right.
- Cell 2**:** An identical circle with a diagonal line, from the top left to bottom right.
- Cell 3**:** A square with an "X" crossing through it.

Second Row:

- Cell 4**:** A combination of a partial square missing the right side, with a triangle pointing upwards.
- Cell 5**:** A partial square missing the left side.
- Cell 6**:** An upward-pointing triangle with a horizontal line intersecting it about one-third from the top.

Third Row:

- Cell 7**:** A mirrored combination of a partial circle on the left and an "L" shaped line on the right.
- Cell 8**:** Similar to Cell 7, a mirrored partial circle on the left with an extended "L" shape on the right.
- Cell 9**:** Contains a question mark, indicating the missing symbol in this cell.

Answer Options

Option A: A combination of a mirrored partial circle on the left with a small "L" shaped line extending from the right side.

Option B: Similar to Option A, but the partial circle has an additional horizontal line.

Option C: A vertical and horizontal series of short parallel lines resembling an incomplete rectangle.

Option D: An "L" shaped line with a mirrored partial circle on the left side, without any additional lines.

The task involves identifying the pattern or logic that governs the symbols in the grid to select the correct option that completes it.] Select the most appropriate option from the four given choices to fill in the question mark, so that a certain pattern is presented:

A: A\nB: B\nC: C\nD: D

IMAGE DESCRIPTION:The image consists of two rows, each containing boxes with geometric shapes, indicating a puzzle pattern.

Top Row:

- First Box**:** Contains a black equilateral triangle centered within a white square.
- Second Box**:** Features two black triangles, one inverted and one upright, forming an hourglass shape, centered inside a pentagon.
- Third Box**:** Shows a black triangle pointing right, inside a horizontal white oval.
- Fourth Box**:** Displays a black rectangle and a larger black semicircle, creating a shape like the letter 'C,' inside a circle.

Bottom Row:

- Fifth Box**:** Contains a black diamond inside a pentagon.
- Sixth Box**:** Contains a question mark, suggesting the missing piece in the pattern.
- Seventh Box**:** This row presents four possible answer choices labeled A to D.
- Option A**:** Depicts a horizontal black oval inside a white heart shape.
- Option B**:** Includes a black shape resembling a curved segment inside a pentagon.
- Option C**:** Features a black diagonal stripe dissecting a white cross.
- Option D**:** Displays an L-shaped black figure inside a pentagon.

Patterns in the top row suggest a sequence or rule, and the task is to determine which option (A, B, C, or D) should replace the question mark based on the continuation of the sequence or pattern.

Figure 14: Part of image caption in LLM experiment.

IMAGE DESCRIPTION:The image presents a visual pattern recognition puzzle consisting of a sequence of geometric shapes organized in two rows.

The objective is to determine the missing shape in the sequence from the given options.

Top Row (Sequence of Shapes)

- First Shape**:** Two identical triangles, one upright and one upside down, overlap each other perfectly to form a star-like shape with a total of six outward-pointing triangular tips.
- The outlines of the triangles intersect, forming an internal hexagon.

- Second Shape**:** Two overlapping triangles, one upright and one upside down, slightly offset and not perfectly aligned, creating a symmetric "X" formation at the center.
- The overlap forms a six-pointed star with more distinct and enhanced intersection lines visible inside.

- Third Shape**:** Two overlapping triangles, one upright and one upside down, arranged symmetrically to create a six-pointed star.
- The internal lines are more connected, forming a central hexagon similar in style to the first shape, with additional internal lines from the triangle overlaps.

- Fourth Shape**:** Two overlapping parallelograms, forming a shape that looks like a hexagon outlined in a bold, overlapping manner.
- The parallelograms intersect to create a pattern with sharp angles and straight lines defining internal quadrilateral shapes.

- Fifth Shape**:** A centralized octagon with an inner pattern of lines creating a star shape, with the corners of the star connecting each side of the octagon from the center.

Missing Shape (Question Mark):** This is the space where the next shape in the sequence would logically fit, represented by a question mark.

Bottom Row (Answer Options)

- Option A**:** Two vertical rectangles with a central pointed triangle merging at the top to form an arrowhead shape.
- The base of the rectangles is flat, with an internal vertical line separating them.

- Option B**:** Two triangles placed side by side with a small gap between them.
- The triangles are upright and identical in size and shape.

- Option C**:** Two parallelograms overlapping each other slightly, tilted sideways.
- The overlapping creates an interesting central intersection with slanting lines.

- Option D**:** An octagon with a central star pattern formed by crossing diagonal lines.
- The star is symmetrical and fits perfectly within the octagon's outline.

Task

The objective is to determine which of the four options (A, B, C, or D) logically completes the sequence in the top row.] From the four given options, choose the most suitable one to fill in the question mark, so that it presents a certain regularity:

A: A\nB: B\nC: C\nD: D

IMAGE DESCRIPTION:The image is a 3x3 grid puzzle with the task of identifying the missing symbol in the bottom right cell. Below the grid are four answer options labeled A, B, C, and D.

Grid Description

Cell 1**: A circle with a diagonal line running from the top left to the bottom right.	Cell 2**: An identical circle with a diagonal line, from the top left to bottom right.	Cell 3**: A square with an "X" crossing through it.
Cell 4**: A combination of a partial square missing the right side, with a triangle pointing upwards.	Cell 5**: A partial square missing the left side.	Cell 6**: An upward-pointing triangle with a horizontal line intersecting it about one-third from the top.
Cell 7**: A mirrored combination of a partial circle on the left and an "L" shaped line on the right.	Cell 8**: Similar to Cell 7, a mirrored partial circle on the left with an extended "L" shape on the right.	Cell 9**: Contains a question mark, indicating the missing symbol in this cell.

Answer Options

- Option A**:** A combination of a mirrored partial circle on the left with a small "L" shaped line extending from the right side.
- Option B**:** Similar to Option A, but the partial circle has an additional horizontal line.
- Option C**:** A vertical and horizontal series of short parallel lines resembling an incomplete rectangle.
- Option D**:** An "L" shaped line with a mirrored partial circle on the left side, without any additional lines.

The task involves identifying the pattern or logic that governs the symbols in the grid to select the correct option that completes it.] Select the most appropriate option from the four given choices to fill in the question mark, so that a certain pattern is presented:

A: A\nB: B\nC: C\nD: D

IMAGE DESCRIPTION:The image consists of two rows, each containing boxes with geometric shapes, indicating a puzzle pattern.

Top Row

- First Box**:** Contains a black equilateral triangle centered within a white square.
- Second Box**:** Features two black triangles, one inverted and one upright, forming an hourglass shape, centered inside a pentagon.
- Third Box**:** Shows a black triangle pointing right, inside a horizontal white oval.
- Fourth Box**:** Displays a black rectangle and a larger black semicircle, creating a shape like the letter 'C,' inside a circle.
- Fifth Box**:** Contains a black diamond inside a pentagon.
- Sixth Box**:** Contains a question mark, suggesting the missing piece in the pattern.

Bottom Row

- Option A**:** Depicts a horizontal black oval inside a white heart shape.
- Option B**:** Includes a black shape resembling a curved segment inside a pentagon.
- Option C**:** Features a black diagonal stripe dissecting a white cross.
- Option D**:** Displays an L-shaped black figure inside a pentagon.

Patterns in the top row suggest a sequence or rule, and the task is to determine which option (A, B, C, or D) should replace the question mark based on the continuation of the sequence or pattern.

Figure 15: Part of image caption in LLM experiment.

Problem

Which of the following four options is the front elevation view of the figure in the question:

GPT-4o

"To solve this problem, we need to determine the front elevation view of the given 3D figure.
Let's break it down:
Identify the visible blocks from the front view: - There are four blocks visible on the front face: D (bottom left), A (middle), and B (top left), and another A on the bottom right.
- The C block is behind the A block in the middle and part of it will not be visible from the front.
Match the arrangement with the options:
- **Option B:** Shows B on the top, A in the middle, and C on the right. This does not match the front view as C is not visible directly from the front.
- **Option C:** Shows B at the top, D at the bottom, and A in the middle. This is incorrect as the positioning does not fit.
- **Option D:** Shows B on the top, A in the middle, and D at the bottom but is missing the extra A block visible from the front.
Verify visibility issue: - From the front, only blocks B, A, and D are visible, and Option A correctly represents this view with the correct alignment.
Therefore, the correct front elevation view is represented by the option that correctly shows the relative positioning and visibility of the blocks.
Answer: **\boxed{A}**"

Qwen2.5-72B-Instruct

To determine the correct front elevation view of the given 3D figure, we need to analyze the structure from the front perspective.
Step 1: Identify the structure from the front.
The topmost layer has one block labeled 'B'.
The middle layer has two blocks labeled 'A' and 'C'.
The bottom layer has three blocks labeled 'D', 'A', and 'B'.
Step 2: Compare with the given options.
- **Option A:** Shows 'B' on top, 'A' in the middle, and 'D' at the bottom. However, it does not show 'C' in the middle layer.
- **Option B:** Shows 'B' on top, 'C' in the middle, and 'D' at the bottom. However, it does not show 'A' in the middle layer.
- **Option C:** Shows 'B' on top, 'A' in the middle, and 'D' at the bottom. This matches the structure from the front.
- **Option D:** Shows 'B' on top, 'A' in the middle, and 'C' at the bottom. However, it does not show 'C' in the middle layer.
Based on the analysis, the correct front elevation view is shown in Option C.
Answer: **\boxed{C}**

Reference Solution

The front view of the figure in the question should be the one with three squares on the left and the letter B on the top, the two faces with letters A and D in the middle, and a blank face on the right, which corresponds to option A. So the correct answer is A

Gemini-2.0-pro-exp-02-05

Step 1: Identify the front face. The front face is the side where blocks A and D are most prominent. Imagine looking directly at the structure from the front. There are two rows of blocks. The top row (front row) visible from the front. Block D is visible at the bottom left. There is another block immediately to the right of D, also visible from the front. So, the bottom row has two blocks. Step 3: Determine the middle layer visible from the front. Block A is visible directly above block D. There is another block immediately to the right of A, also visible from the front. So, the middle row has two blocks. Step 4: Determine the highest layer (top row) visible from the front. Block B is visible directly above block A. There are no other blocks in the top layer visible from the front. So, the top row has one block. Step 5: Summarize the front view of the figure. The front elevation view is: Top row: Block [B] | Middle row: Block [A] | Bottom row: Block [D] |
Step 6: Compare this structure with the given options.
Option A matches this structure exactly: B in the top-left, A below B, D in the middle row, and two blocks in the bottom row.
Option B shows block C instead of A and has an incorrect structure.
Option C has only one block in the middle row.
Option D has only one block in the bottom row.
Answer: **\boxed{A}**

Image Caption

IMAGE DESCRIPTION: The image consists of two main sections. The top part shows a three-dimensional structure of connected cubes, while the bottom part presents four flat arrangements labeled A, B, C, and D for comparison.

Top Section: 3D Cube Structure
- The 3D structure is made up of eight individual smaller cubes.
- **Front View Arrangement:** The structure has a visible front face composed of four cubes. From left to right and top to bottom, these cubes are labeled as B, C, A, and D. The top-left cube is labeled "B". The top-right cube is labeled "C". The bottom-left cube is labeled "D". The bottom-right cube is labeled "A".
- **Middle Row:** Cubes labeled "C" are on the top-left and bottom-left.
- **Bottom Row:** Cubes labeled "D" are on the top-right and bottom-right.
- **Depth and Height:** There is an additional cube directly behind the cube labeled "A", providing depth. Another cube is above the cube labeled "A" on top of cube "B".
Bottom Section: Flat Arrangements
- **Four Options (A, B, C, D):** Each option represents a potential flattened or 2D configuration of the 3D structure.
- Option A: A vertical line of three cubes with "B" at the top, followed by "A", and then "D" at the bottom. A single cube "C" is attached to the right of "A".
- Option B: A vertical line similar to Option A with "B" at the top, followed by "A", and "D" at the bottom. A single cube "C" is attached to the right of "D".
- Option C: A vertical line with "B" on top, followed by "C", then "A", and "D" at the bottom. This sequence is a full vertical line without side attachments.
- Option D: A horizontal line of three cubes starting with "B" on the left. Below "B" is "A", followed by "D" on the same level. A single cube "C" is attached to the right of "A".
- Purpose: The task is to match the 3D structure's visible and implied arrangement with one of the 2D flattened representations labeled A, B, C, or D.

Douba-1.5-Vision-Pro-32k-250115

1. First, understand the concept of the front - elevation view:
- The front - elevation view shows the shape and relative position of the objects when viewed from the front.
2. Analyze the given 3 - D figure:
- Looking at the 3 - D figure from the front, we can see that the top - most cube is B.
- Directly below B, there is a cube A, and to the right of A (at the same level as A) is cube C.
- Below A, there is cube D.
3. Analyze each option:
- Option A: In this option, the relative positions of B, A, C, and D are correct.
- Option B: A is at the top - most position, C is to the right of A at the same level as A, and D is below A.
- Option C: In this option, the position of C is incorrect. C should be at the same level as A but to the right of A, while in this option, C is at the same level as B and to the right of B.
- Option D: In this option, the position of C is incorrect. C should be at the same level as A but to the right of A, while in this option, C is at the same level as B and to the right of B.
Answer: **\boxed{A}**

Claude-3-7-Sonnet-Thinking

<think>Let me analyze this visual reasoning problem carefully.
The top part of the image shows a 3D arrangement of cubes labeled with letters A, B, C, and D. The bottom part shows four different 2D views labeled as options A, B, C, and D.
I need to determine which of these 2D views represent the front elevation view of the 3D arrangement.
First, let me understand the arrangement of cubes in the 3D figure:
- Cube A is positioned in the middle of the arrangement.
- Cube B is positioned on top of A.
- Cube C is positioned to the right of A.
- Cube D is positioned at the lower level, to the left of A.
Now, let's consider what a front elevation view means:
A front elevation view is a 2D representation looking at the object from the front, showing only what would be visible from that perspective.
Let me examine each option:
- Option A: Shows a 2D view with B in the top-left, A in the middle, and D extending to the right.
- Option B: Shows a 2D view with B in the top-left, C in the middle-right, and D in the bottom-left.
- Option C: Shows a 2D view with B in the top-middle, A in the middle-left, and D extending to the right.
- Option D: Shows a 2D view with B in the top-right, A in the middle-right, and D in the bottom-right.
Now, if I look at the 3D figure from the front (assuming the front is where cube D is visible), I would see:
- Cube D at the bottom-left.
- Cube A above D and to the right.
- Cube B above A and to the right.
- Cube C to the right of A.
Therefore, the front elevation view would be Option A.</think>
Front Elevation View Analysis
To solve this problem, I need to identify which 2D view represents the front elevation of the 3D arrangement shown in the top image.
Step 1: Understand the 3D arrangement.
In the 3D figure:
- Cube A is positioned in the middle layer.
- Cube B is positioned on top of A.
- Cube C is positioned to the right of A.
- Cube D is positioned at the lower level, to the left of A.
Step 2: Determine what a front elevation view should show.
A front elevation view shows what you see when looking directly at the object from the front. In this case, I'll assume the front is the face where we can see cube D.
Step 3: Analyze each option.
Looking at the options:
- Option A: B is at the top-left, A is in the middle, and D extends to the right.
- Option B: Shows B at top-left, C is in the middle-right, and D at bottom-left with cubes extending right.
- Option C: Shows B at top-middle, A in the middle-left, and D extending right at the bottom.
- Option D: Shows B at top-right, A in the middle-right, and D at the bottom-right.
Step 4: Match with the 3D arrangement.
When looking at the 3D figure from the front (where D is visible):
- D should be at the bottom-left.
- A should be above and to the right of D.
- B should be above A.
The positioning in Option A matches this arrangement perfectly.
Answer: **\boxed{A}**

Qwen2.5-72B-Instruct

To solve this problem, we need to carefully analyze the 3D structure and compare it to the 2D arrangements provided. Let's break down the steps:
Step 1: Identify the Front View Arrangement.
- The front view of the 3D structure shows cubes labeled B, C, A, and D.
- The arrangement from left to right and top to bottom is: Top row: B (top), C (middle), A (bottom). Bottom row: D (left), A (middle), B (right).
Step 2: Analyze the Depth and Height.
- There is an additional cube directly behind the cube labeled "A".
- Another cube is above the cube labeled "A" on top of cube "B".
Step 3: Compare with the given options.
- **Option A:** Shows B at top-left, A in the middle, and D extending to the right.
- **Option B:** Shows B at top-left, C in the middle-right, and D at bottom-left with cubes extending right.
- **Option C:** Shows B at top-middle, A in the middle-left, and D extending right at the bottom.
- **Option D:** Shows B at top-right, A in the middle-right, and D at the bottom-right.
Step 4: Match with the 3D arrangement.
When looking at the 3D figure from the front (where D is visible):
- D should be at the bottom-left.
- A should be above and to the right of D.
- B should be above A.
The positioning in Option A matches this arrangement perfectly.
Answer: **\boxed{A}**

Figure 16: Solution examples generated by different models. Reference solution and outputs generated by GPT-4o, Qwen2.5VL-72B-Instruct, Gemini-2.0-pro-exp-02-05, Douba-1.5-Vision-Pro-32K and Claude-3.7-sonnet-thinking. Additionally, the image caption and solution from LLMs (Qwen2.5-72B-Instruct) are also illustrated.

Problem

Which of the following four options is the front elevation view of the figure in the question:

A B C D

GPT-4o

To solve the problem, let's analyze the sequence of figures:
 1. **Identify Patterns in Given Figures:**
 - Break down each figure in the sequence to identify patterns or transformations.
 - The first image contains a large square divided into several triangular sections.
 - The second image shows an edge merging and forming a new shape.
 - The third image introduces more fragmented parts.
 - The fourth image repeats a similar pattern with a new twist in fragmentation.
 2. **Look at the Specific Transformations:**
 - Notice possible rotations, reflections, or consistent transformations that occur from one figure to the next.
 - There is a pattern of dividing the square and incrementing triangular or rectangular partitions in each step.
 3. **Predict the Next Transformation:**
 - By examining figures, notice that each set seems to introduce a new line or rotate slightly and changes align with a simple geometric transformation.
 - The consistent theme of geometric division is essential.
 4. **Applying the Observed Pattern:**
 - Look at options A, B, C, and D.
 - Option A features a very intricate pattern with no straightforward changes observed in previous steps.
 - Option B shows a diagonal shape consistent with transformations.
 - Options C and D suggest other structural changes that fit less with the observed pattern.
 5. **Focus on Selection that continues simple progression of incremental shape division or mirror.**
 - Select the Answer.
 6. **Consider all evidence and the goal for minimal transformation matching observed traits.**
Answer: B

Reference Solution

It is found that the closed area of the figure in the question stem is obvious. Considering the number of faces, the number of faces in the figure in the question stem and the option figures is 6, so it is impossible to choose the answer. Considering the refinement of the faces, it is found again that the figure in the question stem basically contains triangular faces. Considering the number of triangular faces, the number of triangular faces in the figure in the question stem is 0, 1, 2, 3, and 4 in sequence, so the number of triangular faces in the figure at "?" should be 5, and only option B meets the requirement. So the correct answer is B.

Doubao-1.5-Vision-Pro-32k-250115

1. First, analyze the characteristics of the figures in the question - set A. Each figure in the question - set is composed of a large square divided into several small right - angled triangles. We can consider the relationship between the lines in the figures, such as the parallel and perpendicular relationships of the line segments.
 - Observe the number of right - angled triangles in each figure: The first figure has 5 right - angled triangles, the second figure has 4 right - angled triangles, the third figure has 3 right - angled triangles, the fourth figure has 2 right - angled triangles, and the fifth figure has 1 right - angled triangle. So, the number of right - angled triangles in the figures in the question - set is decreasing by 1 successively.
 2. Then, analyze the options:
 - Option A: The figure has 2 right - angled triangles, which does not meet the requirement of the number of right - angled triangles decreasing by 1 successively.
 - Option B: The figure has 0 right - angled triangles. According to the rule that the number of right - angled triangles in the figures in the question - set is 5, 4, 3, 2, 1, and then 0, this option meets the rule.
 - Option C: The figure has 4 right - angled triangles, which does not meet the rule.
 - Option D: The figure has 1 right - angled triangle, which does not meet the rule.
Answer: B

Kimi-latest

To solve this problem, we need to identify the pattern in the sequence of the first five figures and then determine which of the options (A, B, C, or D) fits this pattern.
 1. **Observe the Sequence:** Look at the sequence of the first five figures. Notice that each figure is a square divided into four smaller triangles by two diagonals that intersect at the center of the square.
 2. **Identify the Pattern:** Observe how the triangles are shaded in each figure:
 - The first figure has the top-left and bottom-right triangles shaded.
 - The second figure has the top-right and bottom-left triangles shaded.
 - The third figure has the top-left and bottom-left triangles shaded.
 - The fourth figure has the top-right and bottom-right triangles shaded.
 - The fifth figure has the top-left and bottom-right triangles shaded.
 3. **Determine the Pattern:** The pattern alternates between shading the top-left and bottom-right triangles and shading the top-right and bottom-left triangles.
 4. **Apply the Pattern:** Since the fifth figure has the top-left and bottom-right triangles shaded, the next figure in the sequence should have the top-right and bottom-left triangles shaded.
 5. **Match with Options:** Look at the options (A, B, C, D) and identify which one matches the pattern of having the top-right and bottom-left triangles shaded.
 - Option A: Top-left and bottom-right triangles shaded.
 - Option B: Top-right and bottom-left triangles shaded.
 - Option C: Top-left and bottom-right triangles shaded.
 - Option D: Top-right and bottom-left triangles shaded.
 6. **Select the Correct Option:** Both options B and D have the top-right and bottom-left triangles shaded. However, we need to ensure that the pattern is consistent with the previous figures. Since the sequence alternates and the fifth figure has the top-left and bottom-right triangles shaded, the next figure should indeed have the top-right and bottom-left triangles shaded. Both B and D fit this pattern, but we need to choose one. Given the options, we can select either B or D. For the sake of this solution, we will choose option B.
 Therefore, the most suitable option to fill in the question mark is B.

Gemini-2.0-pro-exp-02-05

Analyze the sequence: Observe the five figures provided in the sequence (let's call them F1, F2, F3, F4, F5) and look for a pattern or relationship between them to determine the sixth figure (F6).
 1. **Examine structural features:** Count line segments, analyze shapes, check symmetry, look at intersections, etc. Each figure consistently has 6 internal line segments. The shapes and symmetries vary significantly.
 2. **Look for simple transformations:** Rotations or reflections do not seem to map consecutive images directly (e.g., F1 rotated/reflected is not F2).
 3. **Consider alternating patterns or cycles:** Maybe odd-numbered figures (F1, F3, F5) follow one pattern, and even-numbered figures (F2, F4, F6) follow another.
 4. **Hypothesize a relationship with a fixed offset:** Let's test if there's a consistent transformation T such that $F(n+k) = T(F(n))$ for some k. Let's try $k=3$.
 5. **Check F1 \rightarrow F4:** F1 has a central structure based on connecting midpoints (a diamond shape) plus vertical line segments connecting the diamond to the center point (or near it). Let's denote this structure as $(MD \rightarrow)$. F4 has a central structure based on the two main diagonals (an 'X' shape) plus two smaller diagonal segments ($slope=1$) in the bottom-left and top-right corners relative to the center. Let's denote this structure as $(X \rightarrow \text{shifted})$. The transformation $T1 \rightarrow 4$ changes the base structure ($MD \rightarrow X$) and the added lines ($\rightarrow \text{shifted}$).
 6. **Check F2 \rightarrow F5:** F2 has a complex structure. F5 has one main diagonal ('/ shape) plus three other line segments. Let's denote this as $(/+ L)$. The transformation $T2 \rightarrow 5$ changes the base structure ($F2 \rightarrow /$) and the added lines ($F2 \rightarrow L$). This fits the general idea of "Change Base, Change Additions".
 7. **Predict F3 \rightarrow F6:** F3 has a central structure based on horizontal and vertical midlines ('+ shape) plus two diagonal segments (slope -1) in the top-left and bottom-right corners relative to the center. Let's denote this structure as $(+ \rightarrow \text{shifted})$. Apply the transformation logic: Change Base, Change Additions. Following the pattern seen in $T1 \rightarrow 4$ and potentially $T3 \rightarrow 6$ (where the base also becomes X), let's assume the base (+) changes to (X). What happens to the added lines (diagonals)? In $T1 \rightarrow 4$, changed to /. In $T3 \rightarrow 6$, maybe / changes to ? or (+) changes to (X) and (X) changes to (), the resulting figure F6 would be (X +). This means an 'X' shape from the main diagonals plus vertical line segments. Let's check the options.
 8. **Option A: (X + -) - Diagonal cross plus horizontal connectors.**
 9. **Option B: (/ + complex) - Slash diagonal plus other lines.**
 10. **Option C: (/ + complex) - Slash diagonal plus other lines.**
 11. **Option D: (X + -) - Diagonal cross plus vertical connectors.**
 The predicted structure $(X + -)$ matches Option D exactly.
Conclusion: The pattern $F(n+3) = \text{Transform}(F(n))$ holds consistently, where the transformation involves changing the base structure and the type/position of added lines. Applying this transformation to F3 suggests F6 should be Option D.
Answer: D

Image Caption

IMAGE DESCRIPTION: The image presents a puzzle made up of a series of geometric patterns. The top row consists of five squares in a sequence followed by a question mark indicating a missing pattern in the sixth square. Below this row, there are four options labeled A, B, C, and D, each containing a unique geometric pattern inside a square. **Top Row Pattern Series:**
 - **Option A:** **Square 1** - Contains a large diagonal line from the bottom-left to the top-right, dividing the square into a triangular shape on the bottom-left and - Within the upper-right section, a vertical and horizontal line meet, forming a smaller triangle on a small trapezoid.
 - **Option B:** **Square 2** - Contains a vertical line dividing the square in half with a horizontal line at the mid-top extending partially in. A diagonal runs from the top left to the bottom right corner through these lines, creating several intersecting triangular regions.
 - **Option C:** **Square 3** - Features a central diagonal line from the top right to the bottom left, creating two large triangles in. Additional lines partition these triangles into smaller triangles and trapezoids.
 - **Option D:** **Square 4** - Contains a vertical line dividing the square in half with a horizontal line at the mid-top extending partially in. A diagonal runs from the top left to the bottom right corner through these lines, creating several intersecting triangular regions.
Question Mark Pattern:
 - The pattern has a large diagonal from the top left to the bottom right. An additional horizontal line creates two triangles on either side with a small square or rectangle in the center.
Options:
 - **Option A:** Bottom Row (Options A, B, C, D): Option A is correct.
 - **Option B:** Option B is correct.
 - **Option C:** Option C is correct.
 - **Option D:** Option D is correct.
Task: The task is to determine which of the four options (A, B, C, D) logically continues the sequence established by the first five squares in the top row. Each square and option is enclosed with a solid border, emphasizing the contained geometric pattern. The solution involves recognizing a pattern or rule that connects the top sequence with one of the provided options.

Qwen2.5-72B-Instruct

To solve this visual logical reasoning problem, let's analyze the patterns and transformations occurring from one square to the next in the top row. We will look for consistent changes or rules that can help us predict the next pattern.

Step-by-Step Analysis:

- Square 1 to Square 2:** Initial Pattern. Large diagonal from bottom left to top right. Transformation: The diagonal flips to go from top left to bottom right, and an additional horizontal line is added, creating more segments.
- Square 2 to Square 3:** Initial Pattern. Large diagonal from top left to bottom right with a horizontal line. Transformation: Additional vertical and horizontal lines are added, creating a grid-like structure with a diagonal intersecting it.
- Square 3 to Square 4:** Initial Pattern. Grid-like structure with a diagonal. Transformation: The grid simplifies to a central diagonal from top right to bottom left, with additional lines creating smaller triangles and trapezoids.
- Square 4 to Square 5:** Initial Pattern. Central diagonal from top right to bottom left with additional lines. Transformation: The pattern simplifies further with a vertical line dividing the square in half and a horizontal line at the mid-top, with a diagonal running through them.

Observations:

- Diagonal Flipping:** The direction of the main diagonal alternates between squares.
- Line Addition and Simplification:** Additional lines are added to create more segments, but then the pattern simplifies in the next square.
- Symmetry and Complexity:** The complexity of the pattern increases and then decreases, suggesting a cyclical nature.

Predicting the Next Pattern:

Given the alternating diagonal and the cyclical nature of complexity, the next pattern should:

- Have a diagonal from bottom left to top right (flipping back).
- Be relatively simple, possibly with fewer additional lines.
- Evaluating the Options:
- Option A:** Multiple intersecting diagonals creating a star-like effect. This is too complex and does not follow the simplicity trend.
- Option B:** Large diagonal from top left to bottom right with additional lines forming a mix of shapes. This is a good candidate as it follows the diagonal flip and has a moderate level of complexity.
- Option C:** Intersecting vertical, horizontal, and diagonal lines. This is also complex and does not fit the simplification trend.
- Option D:** Multiple intersecting diagonals creating multiple triangular segments. This is similar to Option A and is too complex.

Conclusion: Option B is the most suitable choice as it follows the diagonal flip and maintains a moderate level of complexity, fitting the observed pattern.

Answer: B

Figure 17: **Solution examples generated by different models.** Reference solution and outputs generated by GPT-4o, Kimi-latest, Gemini-2.0-pro-exp-02-05 and Doubao-1.5-Vision-Pro-32K. Additionally, the image caption and solution from LLMs (Qwen2.5-72B-Instruct) are also illustrated.

Problem

Which one of the following figures is different from the other three:

=	%	≈	÷	!	..	:	?
---	---	---	---	---	----	---	---

A. # B. :: C. + D. ×

Reference Solution

It is found that the number of parts of the graph in the question is: 2, 3, 2, 3, 2, 3, 2. Therefore, the graph should be a three-part graph. A, C, and D are all one-part graphs. Only option B meets the requirement. So the correct answer is B.

Qwen2.5VL-72B

To solve this problem, let's analyze the sequence of symbols and try to identify any patterns or relationships that might help us determine the correct symbol for the eighth box.

Step-by-Step Analysis:

- 1st box: "=" (Equal sign)
- 2nd box: "%" (Percent sign)
- 3rd box: "≈" (Approximately equal sign)
- 4th box: "÷" (Division sign with dots)
- 5th box: "!" (Exclamation mark)
- 6th box: ":" (Therefore symbol)
- 7th box: "(" (Colon)
- 8th box: "?" (Question mark)

Observing Patterns:

The sequence includes several mathematical symbols such as "=", "%", "≈", "÷", "!", ":".

Mathematical Symbols: These symbols are commonly used in mathematical expressions and logic.

Logical Symbols: The ":" (therefore) and "(" (colon) are often used in logical reasoning and notation.

Progression or Relationship: We need to determine if there is a logical progression or relationship between these symbols.

Possible Relationships:

- **Mathematical Operations:** The symbols could represent a sequence of operations or logical steps.
- **Logical Flow:** The ":" (therefore) suggests a conclusion or result, and the "(" (colon) could be used to introduce a condition or explanation.
- **Analyzing the Options:**

 - Option A: "#" (Hashtag or number sign) - This symbol is not typically used in standard mathematical or logical notation.
 - Option B: ":" (Because symbol) - This symbol is used in logical notation to denote "because" or "since".
 - Option C: "+" (Plus sign) - This is a common mathematical operation.
 - Option D: "×" (Multiplication sign) - This is another common mathematical operation.

Conclusion: Given the presence of ":" (therefore) and the use of logical symbols, it seems logical to complete the sequence with a symbol that represents a reason or cause. The symbol ":" (because) fits well in this context, as it complements the ":" (therefore) by providing the reason or cause for the conclusion.

Final Answer: B

Image Caption

The image shows a puzzle with a sequence of symbols arranged in a horizontal row within individual square boxes. Below this sequence, there are four multiple choice options labeled A, B, C, and D, each with a distinct symbol.

Symbol Sequence in Boxes:

- 1. The first box contains an equal sign ("=").
- 2. The second box contains a percent sign ("%").
- 3. The third box contains an approximately equal sign ("≈").
- 4. The fourth box contains a division sign with a dot above and below the slash ("÷").
- 5. The fifth box contains an exclamation mark ("!").
- 6. The sixth box contains a symbol of three dots arranged in a triangular formation, with two dots on the bottom and one on top ("..").
- 7. The seventh box contains a colon ":".
- 8. The eighth box is marked with a question mark ("?").

Multiple Choice Options:

- Option A: "#" (Hashtag or number sign)
- Option B: ":" (Because symbol)
- Option C: "+" (Plus sign)
- Option D: "×" (Multiplication sign)

The task is likely to identify a pattern or relationship among the symbols in the sequence to determine which of the four symbols (A, B, C, or D) belongs in the eighth box marked with a question mark.

GPT-4o

To solve this problem, we need to discern the pattern or sequence in the given symbols. Looking at the sequence of symbols: =, %, ≈, ÷, .., :, ?. Let's analyze each one:

- "=" is a symbol for equality.
- "%" denotes percentage.
- "≈" is the symbol for approximation.
- "÷" represents division.
- ".." signifies factorial.
- ".." implies "therefore."
- ".." generally denotes a ratio or presence of something.

To identify the pattern, consider the sequence types:

- The symbols seem to be mathematical or logical operators.
- Each symbol could represent a series or group classification in math.
- Now, for patterns or relationships:
- Equality, percentage, approximation, division: These are basic mathematical operations or relations.
- Factorial is a mathematical operation with numbers.
- "Therefore" and ratio both have logical or presentational meaning.

Now, let's match options:

- A. "#" - Hash, not quite fitting into the mathematical operation.
- B. ":" - "Because" logically follows "therefore" (quite fitting).
- C. "+" - Addition fits mathematically but not particularly a logical sequence.
- D. "×" - Multiplication, a mathematical operator but not sequence-linked.

Given that ":" means "therefore," the next logical sequence in relationships would be ".." meaning "because," which often appears as a subsequent explanation.

Thus, the pattern aligns more with logical progression from ":" to ":".

Doubao-1.5-Vision-Pro-32k-250115

1. Analyze the nature of the given symbols. The symbols in the sequence "-", "%", "≈", "÷", "..", ":" are all special-purpose symbols used in mathematics, logic, or general notation, not common arithmetic operation symbols like "+" and "×". Among the options: Option C "+" and Option D "×" are basic arithmetic operation symbols for addition and multiplication respectively. They do not fit the pattern of non-arithmetic operation special symbols in the given sequence. Option A "#" is a general symbol often used in programming, tagging, etc., but it has no direct relation to the set of symbols related to mathematical and logical notations in the given sequence. Option B ":" is a logical symbol used to mean "because" in logical and mathematical reasoning, which is in line with the nature of the given set of symbols (such as ":" which means "therefore"). Conclusion: The most suitable symbol to fill in the blank to maintain the pattern is ":" as it belongs to the same category of logical and mathematical notation symbols as the ones in the given sequence.

Gemini-2.0-pro-exp-02-05

Detailed step-by-step derivation:

- Analyze the sequence: "-", "%", "≈", "÷", "..", ":".
- The given sequence consists of seven symbols followed by a question mark: "-", "%", "≈", "÷", "..", ":".
- Identify potential patterns: Look for visual patterns, numerical patterns, or logical relationships between the symbols. One noticeable feature is the presence and number of dots in each symbol.
- Count the dots: Count the number of dots (or small circles) in each symbol: "-" : 0 dots, "%": 2 dots, "≈": 0 dots, "÷": 2 dots, "..": 1 dot, ":" : 3 dots, ":" : 2 dots, ":" : 1 dot.
- Unknown: Unknown.
- Formulate the dot sequence: The sequence of the number of dots is 0, 2, 0, 2, 1, 3, 2, ?.
- Look for patterns within the dot sequence: The sequence 0, 2, 0, 2, 1, 3, 2 doesn't immediately suggest a simple arithmetic or geometric progression. Let's examine alternating positions: Odd positions (1st, 3rd, 5th, 7th): 0, 0, 1, 2. This sequence appears to be non-decreasing (0, 0, 1, 2). Even positions (2nd, 4th, 6th, 8th): 2, 2, 3, 2. This sequence starts with two 2s, then a 3. A plausible continuation would be another 3, forming the pattern (2, 2, 3, 3).
- Determine the required number of dots: Based on the alternating pattern hypothesis, the 8th symbol (an even position) should follow the pattern (2, 2, 3, 3), meaning it should have 3 dots.
- Evaluate the options: Check the number of dots in each of the given options:

 - Option A: "#" : 0 dots
 - Option B: ":" : 3 dots
 - Option C: "+" : 0 dots
 - Option D: "×" : 0 dots

- Select the correct option: Only option B ":" has the required 3 dots.
- Confirm the pattern: If B is correct, the complete dot sequence is 0, 2, 0, 2, 1, 3, 2, 3. The patterns for odd (0, 0, 1, 2) and even (2, 2, 3, 3) positions hold.

o3-mini

To solve the puzzle, let's analyze the sequence of symbols and identify the underlying pattern.

Sequence Analysis:

- "=" (Equal Sign): Represents equality.
- "%" (Percent Sign): Represents a ratio or proportion.
- "≈" (Approximately Equal Sign): Indicates approximation.
- "÷" (Division Sign with Dots): Represents division.
- "!" (Exclamation Mark): Often used to denote factorial in mathematics or emphasis in language.
- ".." (Therefore Symbol): Indicates a conclusion.
- ".." (Colon): Used to introduce a list or explanation.

Identifying the Pattern:

- The sequence transitions from mathematical symbols to logical and punctuation symbols.
- After the "therefore" symbol (..), which signifies a conclusion, the colon (:) is used to introduce an explanation or reason.
- Following this logical progression, the next symbol should logically represent "because," which is denoted by the ":" (Because Symbol).

Evaluating the Options:

- Option A: "#" (Hashtag or number sign, unrelated to the pattern).
- Option B: ":" (Because symbol, fitting the logical continuation).
- Option C: "+" (Plus sign, unrelated to the current logical flow).
- Option D: "×" (Multiplication sign, also unrelated to the current pattern).

Conclusion:

- Option B ":" best fits the established pattern by introducing the "because" symbol, which logically follows the "therefore" symbol and colon.

Figure 18: **Solution examples generated by different models.** Reference solution and outputs generated by GPT-4o, Qwen2.5VL-72B, Gemini-2.0-pro-exp-02-05 and Doubao-1.5-Vision-Pro-32k. Additionally, the image caption and solution from LLMs (o3-mini) are also illustrated.

27

Question:

A B C D

Hint:

To solve this pattern transformation puzzle, observe how the lines in each figure change. In each step, lines slanting to the left are rotating clockwise and increasing in number, while lines slanting to the right are also rotating clockwise but decreasing in number. Apply this pattern to transition from the third figure to the next, focusing on the number and orientation of the lines to determine the correct configuration.

Question:

A B C D

Hint:

Think about the pattern formed by the intersecting shapes in the first set of figures. Notice how the number of sides of the intersection area increases sequentially with each figure. Apply this pattern to the second set of figures to determine the shape of the intersection at the question mark. Which option continues this sequence?

Question:

A B C D

Hint:

Examine how the black squares within the grid move. Focus on the pattern of movement in both the inner and outer parts of the grid. Note that in the inner 3x3 grid, except for the center square which remains fixed, the other two black squares shift positions in a specific direction and sequence. In the outer grid, observe how the black squares shift consistently in a particular direction. Which option continues this pattern?

Question:

A B C D

Hint:

Consider how the elements of each figure are composed. Observe that in the sequence on the left, the third figure is a combination of parts from the first two figures. Specifically, the bottom half of the third figure is identical to the bottom half of the first figure, and the top half is identical to the top half of the second figure. Apply this pattern to the sequence on the right. Identify the parts of figures on the right first.

Question:

A B C D

Hint:

Focus on counting specific shapes within each figure. Pay attention to how many of a particular geometric shape appears consistently across the options. Remember, some elements like shading may be designed to distract you, so concentrate on the count of the shapes themselves.

Question:

A B C D

Hint:

Consider the problem as one involving grouping and classification. Focus on the characteristics of the shapes provided. Notice that all shapes are divided into distinct closed regions. While the number of these regions varies, a more subtle criterion is at play. Pay attention to the largest region in each shape. Examine its symmetry: does it exhibit rotational or reflective symmetry?

Question:

<td><td><td></td></td></td>	<td><td></td></td>	<td></td>	
1	2	3	4
<td><td></td><td></td></td>	<td></td> <td></td>		
5	6		

A: ①④⑥, ②③⑤
B: ①⑤⑥, ②③④
C: ①②③, ④⑤⑥
D: ①②④, ③⑤⑥

Hint:

Among the following options, which one conforms to the pattern change of the given figure is:

Question:

A B C D

Hint:

Observe the given figures and pay attention to their symmetry properties. The six figures in the title are all symmetrical figures, but their symmetry types are different. Try to divide the figures into two categories: centrosymmetry and axisymmetry to help identify the regularity between figures. Through this classification, we can find the attribution of each figure and get the correct answer.

Question:

A B C D

Hint:

Consider the arrangement and symmetry of the elements. Focus on the properties of the shapes and their composition. In the first group, the shapes exhibit axial symmetry and are composed of letters. Analyze the letters in both groups, noting the symmetry and repetition. When examining the options, identify the one that maintains the same symmetrical properties and letter composition as the examples given.

Figure 19: Examples of hint prompts. Hint prompts are provided to guide reasoning without revealing the final answer directly.

Table 4: Hyperparameter Settings for SFT Training Stage.

	Qwen2.5-VL-7B-Instruct-SFT	InternVL2.5-38B-SFT
pretrain model	Qwen2.5-VL-7B-Instruct	InternVL2.5-38B
learning rate	0.5e-5	2e-5
batch size	64	128
optimizer	AdamW	AdamW
lr scheduler	cosine	cosine
image strategy	image_max_pixels=262144	max_dynamic_patch=6
warmup ratio	0.1	0.03
max epochs	1	1
bf16	True	True

Table 5: Hyperparameter Settings for RL Training Stage.

	Qwen2.5-VL-7B-Instruct-RL	InternVL2.5-38B-RL
pretrain model	Qwen2.5-VL-7B-Instruct	InternVL2.5-38B
RL Algorithm	rloo	rloo
batch size	128	64
rollout batch_size	256	128
temperature	1	1
n samples per prompt	16	8
prompt max len	1024	4096
generate max len	3000	3000
bf16	True	True
actor learning rate	1e-6	1e-6
init kl coef	0	0

A	B	C	D		

From the four given options, select the most suitable one to fill in the place of the question mark, so that it presents a certain regularity:

QwenVL2.5-7B-Instruct

To solve this problem, we need to identify the pattern in the given sequence of circles and apply it to find the correct option for the question mark.
Let's analyze the sequence step by step:
1. **First Row:** - The first circle has a triangle pointing left.
- The second circle has a triangle pointing down.
- The third circle has a triangle pointing right.
2. **Second Row:** - The first circle has a square pointing up.
- The second circle has a square pointing down.
- We need to determine the orientation of the square in the third circle.
3. **Pattern Recognition:** - In the first row, the triangles are pointing in a cyclic order (left, down, right).
- In the second row, the squares are also pointing in a cyclic order (up, down).
4. **Applying the Pattern:** - If the squares continue the same cyclic pattern as the triangles, the next direction after down should be to the right.
5. **Checking the Options:** - Option A: The square is pointing up.
- Option B: The square is pointing down.
- Option C: The square is pointing right.
- Option D: The square is pointing left.
Since the pattern suggests the square should point to the right, the correct option is: **Answer: C**

QwenVL2.5-7B-Instruct-RL

<think> To solve this problem, we need to identify the pattern in the sequence of shapes. Let's analyze the given images step by step.
1. The first row of images shows a sequence of circles with triangles inside them.
- The first circle has a triangle pointing left.
- The second circle has a triangle pointing down.
- The third circle has a triangle pointing right.
The pattern of the triangles is alternating directions (left, down, right).
2. The second row of images shows a sequence of circles with squares inside them.
- The first circle has a square on the top right corner.
- The second circle has a square on the bottom right corner.
- The third circle has a square on the top left corner.
The pattern of the squares is alternating positions (top right, bottom right, top left).
3. We need to determine the position of the square in the third circle of the second row.
- The pattern alternates between top right, bottom right, and top left.
- The previous positions were top right and bottom right.
- Therefore, the next position should be top left.
Now, let's look at the options provided:
- A: A (This option does not match the pattern of the squares.)
- B: B (This option does not match the pattern of the squares.)
- C: C (This option does not match the pattern of the squares.)
- D: D (This option matches the pattern of the squares, as the square is on the top left corner.)
Thus, the most suitable option to fill in the question mark is D.

Figure 20: Comparison of model outputs before and after RL training stage for Qwen2.5-VL-7B.


A

B

C

D

Among the options below, which one can be obtained by flipping the given figure vertically:

QwenVL2.5-7B-Instruct
X

To determine which option can be obtained by flipping the given figure vertically, we need to follow these steps:
 1. Identify the original figure.
 2. Visualize the vertical flip of the original figure.
 3. Compare the flipped figure with the given options.

Let's start with the original figure:
 The original figure is shown in the leftmost image (labeled as "D"). It has the letters "D", "A", and "K" arranged in a specific pattern on the surface of a dodecahedron.

When we flip this figure vertically, the arrangement of the letters will change. Specifically, the letter "D" will move to the bottom, the letter "A" will move to the top, and the letter "K" will move to the right side.

Now, let's compare this with the given options:

- Option A: The letters are arranged as "K", "V", "D".
- Option B: The letters are arranged as "D", "V", "K".
- Option C: The letters are arranged as "D", "V", "K".
- Option D: The letters are arranged as "K", "V", "D".

From the description of the vertical flip, we see that the correct arrangement after the flip should be "D", "V", "K". This matches with Option B.

Therefore, the correct answer is: B

QwenVL2.5-7B-Instruct
✓

<think> To determine which option can be obtained by flipping the given figure vertically, we need to analyze the spatial relationships and the orientation of the letters within the figure.

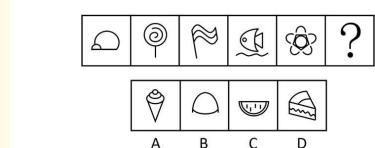
1. **Identify the Original Figure**: The original figure has the letters "D", "K", and "V" arranged in a specific pattern on the surface of a polyhedron. The letter "D" is on the left side, "K" is on the right side, and "V" is in the middle.
2. **Understand Vertical Flipping**: When we flip the figure vertically, the top part of the figure will become the bottom part, and the bottom part will become the top part. The letters will also change their positions relative to each other.
3. **Analyze Each Option**:
 - **Option A**: In this option, the letter "D" is on the left side, "K" is on the right side, and "V" is in the middle. This is the same as the original figure, but it is not a vertical flip.
 - **Option B**: In this option, the letter "D" is on the left side, "K" is on the right side, and "V" is in the middle. This is the same as the original figure, but it is not a vertical flip.
 - **Option C**: In this option, the letter "D" is on the left side, "K" is on the right side, and "V" is in the middle. This is the same as the original figure, but it is not a vertical flip.
 - **Option D**: In this option, the letter "D" is on the left side, "K" is on the right side, and "V" is in the middle. This is the same as the original figure, but it is a vertical flip.
4. **Conclusion**: The only option that shows the figure after a vertical flip is Option D.

</think>

<answer> D </answer>

Figure 21: Comparison of model outputs before and after RL training stage for Qwen2.5-VL-7B.

31



Choose the most appropriate option from the given four options to fill in the question mark, so that it presents a certain regularity:

QwenVL2.5-7B-Instruct 

To solve this problem, we need to identify the pattern or rule that governs the sequence of images in the top row and then apply that same rule to find the correct image for the question mark.
Let's analyze the sequence of images in the top row:
1. The first image is a half-circle.
2. The second image is a lollipop.
3. The third image is a flag.
4. The fourth image is a fish.
5. The fifth image is an atom symbol.
Now, let's look at the images in the bottom row:
A: An ice cream cone.
B: A half-circle.
C: A slice of watermelon.
D: A slice of cake.
We need to determine which of these images fits the pattern established by the top row. Let's consider the shapes and their transformations:
The first image (half-circle) is simple.
The second image (lollipop) has a spiral shape.
The third image (flag) has a wavy line.
The fourth image (fish) has a curved body.
The fifth image (atom symbol) has a complex structure with multiple parts.
The pattern seems to be the transformation of basic shapes into more complex forms. The half-circle becomes a lollipop, which has a spiral shape. The flag has a wavy line, and the fish has a curved body. The atom symbol is complex with multiple parts.
Looking at the bottom row, the image that fits this pattern is the one that transforms a simple shape into a more complex form. The half-circle (image B) matches this description as it is the simplest shape in the top row and corresponds to the simplest shape in the bottom row.
Therefore, the correct answer is:
Answer: B

QwenVL2.5-7B-Instruct 

<think> To solve this problem, we need to identify the pattern in the sequence of images. Let's analyze the images step by step:
1. The first image is a half-circle.
2. The second image is a lollipop.
3. The third image is a flag.
4. The fourth image is a fish.
5. The fifth image is a flower.
6. The sixth image is a question mark, which we need to determine.
Now, let's look at the options provided:
A: An ice cream cone.
B: A half-circle.
C: A watermelon slice.
D: A piece of cake.
We need to identify which of these options fits the pattern established by the previous images. Let's analyze the pattern:
The first image is a half-circle.
The second image is a lollipop, which has a circular shape but with a stick.
The third image is a flag, which has a distinct shape with vertical and horizontal sections.
The fourth image is a fish, which has a distinct shape with a tail and fins.
The fifth image is a flower, which has a distinct shape with petals and a center.
The pattern seems to be that each image has a distinct and recognizable shape. Now, let's look at the options:
Option A (ice cream cone) has a distinct shape with a cone and a scoop.
Option B (half-circle) is already seen in the first image.
Option C (watermelon slice) has a distinct shape with a triangular section.
Option D (piece of cake) has a distinct shape with a triangular section.
Since the pattern is that each image has a distinct and recognizable shape, and the options A, C, and D all have distinct shapes, we need to choose the one that fits the pattern consistently. The most consistent and distinct shape among the options is the ice cream cone.
Therefore, the most appropriate option to fill in the question mark is:
A

Figure 22: Comparison of model outputs before and after RL training stage for Qwen2.5-VL-7B.

				?
--	--	--	--	---

From the four given options, choose the most suitable one to fill in the question mark, so that it presents a certain pattern:

InternVL2_5-38B

X

To solve this visual logical reasoning problem, let's analyze the patterns in the series.

1. Observing the first three shapes:

- The first shape shows a segment (clockwise) from the upper left position going left.
- The second shape shows a segment from the top center going left.
- The third shape shows segments from the bottom left and top left going left.

2. The fourth shape shows a segment from the upper left position going down.

Now, let's identify the progression:

- The segments appear to be moving around the shape in a clockwise manner.
- Each subsequent shape adds an additional segment or changes the position of existing segments.

Given this sequence and the changes between shapes:

- The next in the series should logically follow the pattern of adding or changing segments around the shape in a clockwise direction.

Inspecting the options:

- Option A shows a shape where segments are from the upper left, with one part of the right bottom, which looks like a continuation of the previous pattern.

Hence, option A fits the pattern best.

Answer: **A**

InternVL2_5-38B-RL

✓

<think>

To solve this problem, we need to identify the pattern in the given sequence of shapes and apply it to the next shape.

- The first shape has a line segment positioned at the top-right corner of the circle.
- The second shape has the line segment rotated 90 degrees counterclockwise to the top-left corner.
- The third shape has the line segment rotated 90 degrees counterclockwise to the bottom-left corner.
- The fourth shape has the line segment rotated 90 degrees counterclockwise to the bottom-right corner.

The pattern here is a 90-degree counterclockwise rotation of the line segment in each subsequent shape.

Now, let's examine the options:

- Option A: The line segment is at the top-right corner.
- Option B: The line segment is at the top-left corner.
- Option C: The line segment is at the bottom-left corner.
- Option D: The line segment is at the bottom-right corner.

Since the pattern involves a 90-degree counterclockwise rotation, the next position should be the top-right corner, which corresponds to Option A.

Therefore, the correct answer is D.

</think>

Figure 23: Comparison of model outputs before and after RL training stage for InternVL-2.5-38B.

Divide the following six figures into two categories, ensuring that each category of figures has its own common characteristics or patterns. The correctly classified option is:

InternVL2_5-38B
✗

To determine the correct classification of the given figures, let's analyze the patterns and characteristics of each figure:

- 1. **Figure ①**: A four-sided polygon with a concave shape on the lower part having five vertices if you include internal points where it turns.
- 2. **Figure ②**: Another four-sided figure but with a concave shape on two sides, making it pentagon-like as it consists of five vertices.
- 3. **Figure ③**: A four-sided polygon with a concave part leading to five vertices total when counting internal bends.
- 4. **Figure ④**: A six-sided polygon (hexagon-like) shape with two concave sections leading to seven total vertices.
- 5. **Figure ⑤**: Similar to Figure ②, a four-sided figure with concave sides summing to five vertices.
- 6. **Figure ⑥**: A four-sided polygon but with one concave section and one protrusion, also totaling five vertices.

Using the characteristics:

- **Vertices Count**: ①, ③, ⑤, ⑥ have five vertices. ②, ④ have five vertices as well, but they look different.
- **Pattern of Concave and Convex sections**: Figures ④, ⑤, and ⑥ are all four-sided (quadrilateral) and have similar concave/convex patterns that slightly complicate their counts but still make them consistent in shape complexity.

With this consideration in mind, we can classify based on quadrilateral -like characteristics with specific concave/convex patterns:

- Four-sided shapes with somewhat complex patterns (④, ⑤, ⑥),
- Shapes differing in vertex count and distinct concave sections (①, ②, ③).

Thus the classification is:

Option B: ①③⑥, ②④⑤

Answer: \boxed{B}

InternVL2_5-38B-RL
✓

<think>

To solve this problem, we need to identify common characteristics or patterns among the figures.

1. Figure 1: A shape with an "L" turn and a curve.
2. Figure 2: A shape with a curve and an "S" turn.
3. Figure 3: A shape with an "L" turn and a straight line.
4. Figure 4: A shape with an "L" turn and a curve.
5. Figure 5: A shape with an "L" turn and a straight line.
6. Figure 6: A shape with a curve and an "S" turn.

We can observe the following:

- Figures 1, 3, and 4 have an "L" turn and a curve.
- Figures 2, 5, and 6 have a curve and an "S" turn.

Therefore, the common characteristic for the first category (①③④) is the "L" turn and a curve, and for the second category (②⑤⑥) is the curve and an "S" turn.

</think>

<answer> A </answer>

Figure 24: Comparison of model outputs before and after RL training stage for InternVL-2.5-38B.

34

X

Which of the following options conforms to the given pattern of change in the figures:

InternVL2_5-38B

To determine which option conforms to the given pattern of change in the figures, let's analyze the sequence step by step.

- First Figure:** The figure contains a square cut out from it.
- Second Figure:** The square cut out appears to be rotated slightly clockwise.
- Third Figure:** The square cut out is rotated slightly more clockwise.
- Fourth Figure:** The square cut out is rotated again, continuing the clockwise rotation pattern each time.
- Fifth Figure:** Following the established pattern, the square cut out should continue to rotate clockwise to the next logical position.

Now, let's evaluate each of the options:

- Option A:** This figure retains the unrotated square cut out. Therefore, it does not follow the pattern.
- Option B:** The square cut out rotates the same way the previous figures are following, with additional rotations on each step. This seems to be correct for continued rotation pattern.
- Option C:** This figure's square cut out matches no existing rotations pattern, and follows a different path of rotation.
- Option D:** The square cut out is rotated very different from the sequence.

By verifying the patterns provided and comparing each option, Option B fits the rotation pattern exactly as continued clockwise.

Answer: **B**

✓

<think>

To solve this problem, we need to identify the pattern of change in the given sequence of figures.

- The first figure shows an octagon with a right angle inside it.
- The second figure shows the same octagon with the right angle rotated 45 degrees clockwise.
- The third figure shows the right angle rotated another 45 degrees clockwise.
- The fourth figure shows the right angle rotated another 45 degrees clockwise.
- The fifth figure shows the right angle rotated another 45 degrees clockwise.

From the above observations, we can see that the right angle inside the octagon is rotated 45 degrees clockwise in each subsequent figure.

Now, we need to identify which of the given options matches this pattern.

Option A: The right angle is in the same position as the first figure.

Option B: The right angle is rotated 45 degrees clockwise from the first figure.

Option C: The right angle is rotated 90 degrees clockwise from the first figure.

Option D: The right angle is rotated 135 degrees clockwise from the first figure.

Since the pattern involves a 45-degree clockwise rotation in each step, the next figure should have the right angle rotated 135 degrees clockwise from the first figure.

Therefore, the correct option is D.

</think>

Figure 25: Comparison of model outputs before and after RL training stage for InternVL-2.5-38B.