# MATHGLANCE: Multimodal Large Language Models Do Not Know Where to Look in Mathematical Diagrams

Yanpeng Sun[*1], Shan Zhang[*†2], Wei Tang[3], Aotian Chen[4]
Piotr Koniusz[5], Kai Zou[6], Yuan Xue[‡4], Anton van den Hengel[‡2]
[1]National University of Singapore, [2]Australian Institute for Machine Learning
[3]Nanjing University of Science and Technology, [4] Ohio State University, [5]Data61♥CSIRO, [6]NetMind.ai
[*]Core contribution    [†]Project lead    [‡] Corresponding author
Project Page: https://mathglance.github.io/

## Abstract

*Diagrams serve as a fundamental form of visual language, representing complex concepts and their interrelationships through structured symbols, shapes, and spatial arrangements. Unlike natural images, their inherently symbolic and abstract nature poses significant challenges for Multimodal Large Language Models (MLLMs). However, current benchmarks conflate perceptual and reasoning tasks, making it difficult to assess whether MLLMs genuinely understand mathematical diagrams beyond superficial pattern recognition. To address this gap, we introduce MATHGLANCE, a benchmark specifically designed to isolate and evaluate mathematical perception in MLLMs. MATHGLANCE comprises 1.2K images and 1.6K carefully curated questions spanning four perception tasks: shape classification, object counting, relationship identification, and object grounding, covering diverse domains including plane geometry, solid geometry, and graphical representations. Our evaluation of MLLMs reveals that their ability to understand diagrams is notably limited, particularly in fine-grained grounding tasks. In response, we construct GeoPeP, a perception-oriented dataset of 200K structured geometry image-text pairs explicitly annotated with geometric primitives and precise spatial relationships. Training MLLM on GeoPeP leads to significant gains in perceptual accuracy, which in turn substantially improves mathematical reasoning. Our benchmark and dataset establish critical standards for evaluating and advancing multimodal mathematical understanding, providing valuable resources and insights to foster future MLLM research.*

## 1. Introduction

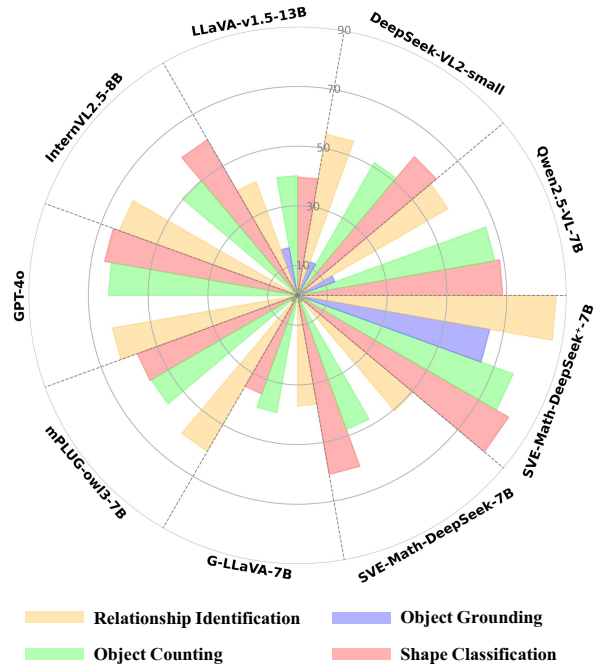Computer vision research has traditionally focused on developing models that can perform fundamental perception



Figure 1. Performance comparison of current Multimodal Large Language Models (MLLMs) on MATHGLANCE reveals a stark contrast: while humans can solve fundamental perception tasks "at a glance", these tasks remain highly challenging for MLLMs—particularly in fine-grained grounding. SVE-Math-DeepSeek[+]-7B, trained with the proposed perception-oriented GeoPeP, shows significantly enhanced perceptual capabilities.

tasks [17, 30] such as object detection [36, 56], segmentation [28, 37], and spatial reasoning [12]. These capabilities form the foundation for higher-level visual reasoning and decision-making. However, recent advances in Multimodal Large Language Models (MLLMs) [1, 32, 42] have shifted the paradigm toward vision-language integration, where visual information is processed through natural language interfaces rather than specialized perception modules.
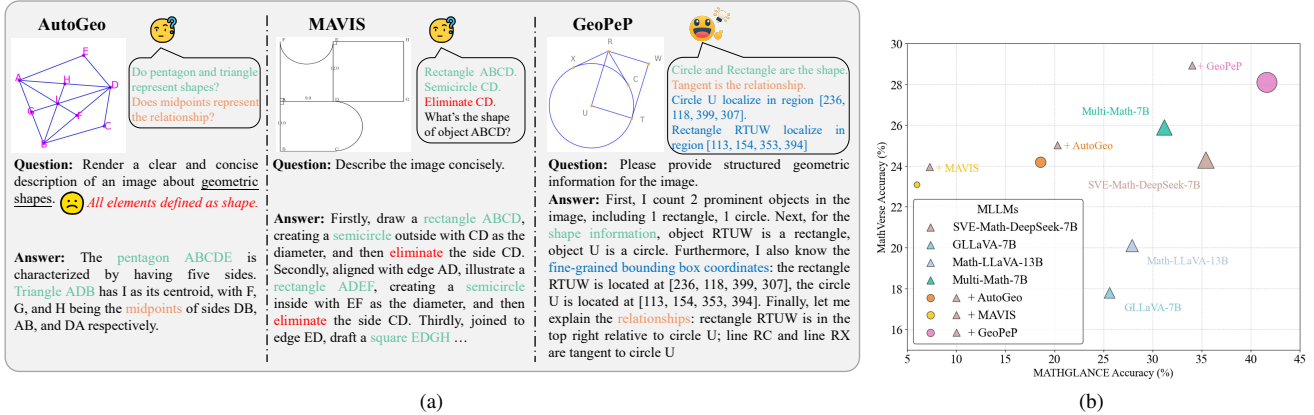
1

Figure 2. Illustration of diagram-caption-based alignment training datasets, w.r.t., AutoGeo, MAVIS, and our proposed GeoPeP (Fig. 2a). Fig. 2b demonstrates a positive correlation between low-level perception and high-level reasoning tasks, evaluated on MATHGLANCE and MathVista. Clear diagram perception leads to substantial improvements in mathematical reasoning performance.

Despite their remarkable successes in general vision tasks, current MLLMs face considerable challenges in accurately interpreting mathematical diagrams—abstract visual representations characterized by precise geometric structures and symbolic notations [34, 53]. These diagrams provide a unique perceptual challenge: the explicit structural constraints require genuine visual interpretation rather than pattern recognition or textual heuristics, positioning them as an essential benchmark for evaluating the true visual perception capabilities of multimodal models. Such diagrams are pervasive across educational contexts, scientific communication, and problem-solving in STEM disciplines. Developing models capable of understanding symbolic information in diagrams, is a critical milestone in advancing machine intelligence [13, 14].

Recent benchmarks like MathVista [34] and MathVerse [53] have attempted to evaluate mathematical visual reasoning in MLLMs. However, these works conflate perception with reasoning by assessment on complex reasoning tasks that combine diagram interpretation with numerical computation and proof generation. Thus, it remains unclear whether the performance truly reflects the models' ability to comprehensively understand the symbolic information in diagrams. Misinterpretations at the perception level propagate errors downstream, resulting in incorrect reasoning outcomes and frequent model hallucinations [7, 23, 44].

To address this critical gap, we design a benchmark to isolate and rigorously evaluate mathematical perception in MLLMs. Our MATHGLANCE features problems that humans can solve "at a glance" without extensive reasoning. The benchmark contains 1,198 images and 1,609 carefully designed perception-oriented questions across four distinct task categories: shape classification, object counting, relationship identification, and object grounding. These tasks span three core mathematical domains—plane geometry, solid geometry, and graphical representations (line, bar,

and pie graphs)—and question formats (multiple-choice, true/false, and free-form).

Through comprehensive evaluation of current MLLMs on MATHGLANCE, we explore the following questions:

1. ***Do current MLLMs genuinely perceive mathematical diagrams?*** Fig. 1 presents the performance of eight current MLLMs—comprising six generic models (both open-source and proprietary) and two math-specific model—across four perception-focused tasks in MATHGLANCE (see Tab. 2 for more results). Among the generic MLLMs, Qwen2.5-VL-7B [6] achieves the highest average accuracy, followed by GPT-4o. But, those models exhibit severe hallucination issues, often responding to simple questions with unnecessarily long chain-of-thought (CoT) reasoning and generating irrelevant image content (see Fig.6). The math-specific SVE-Math-DeepSeek [55], which incorporates a domain-specific visual encoder trained on multi-task geometric primitive detection, achieves comparable performance to Qwen2.5-VL-7B on shape classification and relationship identification tasks. Nevertheless, all evaluated models struggle significantly with fine-grained bounding box grounding. Even those trained on large-scale grounding annotations with nearly 1T natural images—such as Qwen2.5-VL-7B and DeepSeek-VL2-Small—achieve less than 20% accuracy on this task.

To alleviate this issue, we propose GeoPeP, a **Geo**matic **PerceP**tion-oriented dataset. Unlike existing image-caption-based visual training datasets (*e.g.*, MAVIS [54], AutoGeo [21]), GeoPeP provides explicit structured diagram information, including shapes, attributes, locations, and relationships (see Fig.2a). We first train only the vision-language projector of SVE-Math-DeepSeek using this alignment dataset. Then, we construct instruction conversational data for a second self-fine-tuning (SFT) stage, jointly updating both the projector and the
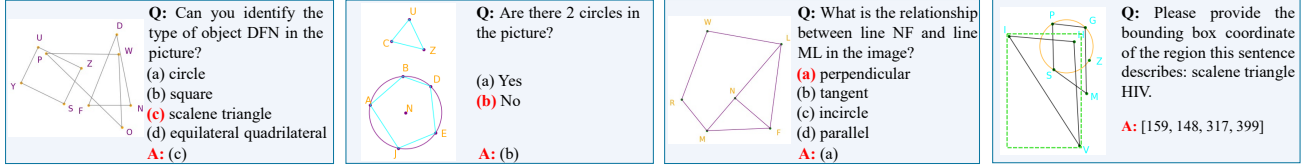
Figure 3. Sampled MATHGLANCE examples from plane geometry w.r.t. each question-answer (Q&A) type, covering shape classification, object counting, relationship identification, and object grounding (from left to right). The green dotted bounding boxe is shown for illustration purposes only and are not provided as input to the models.

large language model. The resulting model, SVE-Math-DeepSeek $^+$, achieves the highest performance among both generic and math-specific models, particularly excelling in grounding tasks with a +79% improvement.

2. ***Does stronger perceptual ability lead to better reasoning performance?*** We hypothesize that the more a model can "see", the more it can "reason" —that is, clear diagram perception would support stronger reasoning ability. As shown in Fig.2b, training with our proposed alignment dataset significantly improves perception performance compared to AutoGeo and MAVIS (by 23.0% and 35.6%, respectively). In contrast, the other two variants perform worse than the baseline (SVE-Math-DeepSeek-7B), particularly with MAVIS diagram-caption pairs, due to high uncertainty and out-of-distribution nature of real-world geometric diagrams. For a fair comparison, we use the same reasoning training dataset (MathV360K [41]) and apply our instruction-style perception sampels to further fine-tune each model. Our model achieves 28.1% accuracy on the MathVista reasoning benchmark—an improvement of +~4% over others. This gain is non-trivial, considering that Multi-Math-7B [38], trained on 600K samples (300K image captions+300K CoT rationales) with additional reinforcement learning, achieves only 26.9% accuracy. Our model accurately identifies relevant visual elements, enabling it to generate more valid and faithful reasoning steps (see model responses in Figs. 7/15-16).

In summary, our contributions are as follows:

- We introduce MATHGLANCE, a comprehensive benchmark for evaluating mathematical perception across geometry and graphical domains. MATHGLANCE isolates perception capabilities through carefully designed tasks and the quality of the data is ensured via rigorous human verification.
- We provide the systematic investigation of how perception capabilities influence mathematical reasoning, quantitatively analyzing MLLMs' perceptual sensitivity to various factors and establishing the correlation between perception accuracy and reasoning performance (see §3.3).
- We construct GeoPeP, a **Geo**matic **P**erce**P**tion-oriented dataset comprising 200K high-quality geometry image-text pairs (100K diagram-caption + 100K conversation) that explicitly struct diagrams into shapes, attributes, lo-

cations and relationships. GeoPeP provides precise visual grounding cues and minimizes ambiguities common in existing caption-based datasets, which can significantly improve MLLMs' visual perception accuracy and downstream reasoning performance.

We envision MATHGLANCE as a valuable testbed to bridge the gap between low-level visual perception and high-level mathematical reasoning, ultimately facilitating the development of more robust and reliable multimodal models for mathematical understanding.

## 2. MATHGLANCE

This section first introduces MATHGLANCE (§2.1) with a data distribution analysis. We then describe the synthetic construction process for plane geometry, as well as the reformatted CLEVR [24] and FigureQA [25] datasets for solid geometry and graphs in §2.2. Finally, we introduce the perception-oriented GeoPeP, and highlight its unique characteristics in comparison to existing mathematical visual alignment training datasets (§2.3).

### 2.1. Overview

Existing mathematical visual reasoning benchmarks often conflate perception with higher-level reasoning tasks, such as numeric calculations and proof generation. We introduce MATHGLANCE, a novel benchmark designed to evaluate MLLMs' mathematical perception-demanding abilities through both quantitative and qualitative analysis across coarse-to-fine granularity levels.

**Data Composition.** Detailed statistics for data composition are presented in Tab. 1. To comprehensively assess the perceptual abilities of MLLMs in mathematical contexts, our evaluation images cover plane geometry (66%), solid geometry (20%), and graphs (14%), including lines, bars, and pie charts. To facilitate MLLM evaluation, we formulate all tasks—except for bounding box grounding (free-form)—as multiple-choice or true/false question-answering problems. In total, we contribute 1,609 questions and 1,198 unique images, ensuring an even distribution across the different perception tasks.

**Categorization.** The benchmark encompasses *shape classification*, *object counting*, *relationship identification*, and *object grounding*. Fig. 3 shows example illustrations of

Table 1. Key statistics of MATHGLANCE are summarized in Tab. 1a, and the subject-task distribution is illustrated in Fig. 1b.

| Statistic | Number |
|---|---|
| Total questions | 1,609 |
| - Multiple-choice questions | 893 (55.5%) |
| - Free-form questions | 406 (25.2%) |
| - True-false questions | 310 (19.3%) |
| Unique number of images | 1,198 |
| Unique number of questions | 1,514 |
| Unique number of answers | 380 |
| Total questions | 1609 |
| - Classification questions | 489 (30.4%) |
| - Counting questions | 401 (24.9%) |
| - Relation questions | 313 (19.5%) |
| - Grounding questions | 406 (25.2%) |
| Maximum question length | 200 |
| Maximum answer length | 4 |
| Maximum choice number | 4 |
| Average question length | 118.0 |
| Average answer length | 1.8 |
| Average choice number | 3.5 |
| Average question length | 16.09 |
| Average answer length | 1.21 |
| Average choice number | 3.40 |

(a)



(b)

plane geometry (see §B for additional examples). Key features are as follows:

- Shape classification is a classic vision task where the model identifies an object's class based on its attributes, *i.e.*, vertices, material, color, and size. Our dataset includes a diverse set of geometric categories, comprising 16 basic shapes for plane geometry, 3 CLEVR-defined [24] objects for solid geometry, and 5 graphical elements as defined in FigureQA [25].
- Object counting requires models to determine either the total number of objects in an image or a specific shape count, *i.e.*, the number of circles or triangles present.
- Relationship identification evaluates models' understanding of 4 spatial and over 10 mathematical relationships between pairs of geometric primitives.
- Object grounding evaluates fine-grained localization by requiring MLLMs to accurately predict the top-left and bottom-right coordinates in the format (x1, y1, x2, y2) for an object within the image. This ensures that models can precisely identify and localize geometric structures based on textual descriptions.

These four tasks are designed to isolate and assess fundamental perception capabilities, distinguishing low-level visual recognition from high-level reasoning.

## 2.2. Dataset Construction

### 2.2.1. Synthetic Data Engine for Plane Geometry

We design a data generation engine to synthesize real-world geometric images along with structured JSON annotations. This pipeline enables controlled shape generation, relation-

ship modeling, and visual attribute assignment, ensuring diverse and well-balanced datasets across perception tasks.

**Structured Annotations for Geometric Primitives.** Fig. 4 describes the entire generation process. Inspired by AlphaGeometry [43], we use geometric clauses as fundamental units to construct complex plan geometric figures. A geometric clause is a formalized description of basic geometric objects and mathematical relationships, along with their properties or attributes, a.k.a. prerequisite points. We first construct two geometry substrate pools, one containing 16 different geometric shapes (*i.e.*, scalene triangle, isosceles triangle, square, rectangle, parallelogram, isosceles trapezoid, right trapezoid, pentagon, circle . . . ellipse and segment), and the other defining 10 mathematical relationships (*i.e.*, on, intersection, parallel, perpendicular, tangent, incircle, . . . reflection and angle bisector). We then randomly sample one or more substrates from these pools and pass them through a verifier, which ensures that logically paired shapes and relationships are preserved in the construction of valid geometric images. The verifier makes decisions based on either manually designed rules, fundamental mathematical knowledge, or prerequisite points. For example, parallel lines cannot intersect, making an intersection clause invalid. Similarly, angle trisection requires a predefined angle, which cannot be constructed without enough prerequisite points for angle consturction. In general, the chosen relationships are enforced by introducing additional shapes into the figure, ensuring that the specified relationships are accurately maintained and geometrically consistent throughout the construction process. We finally
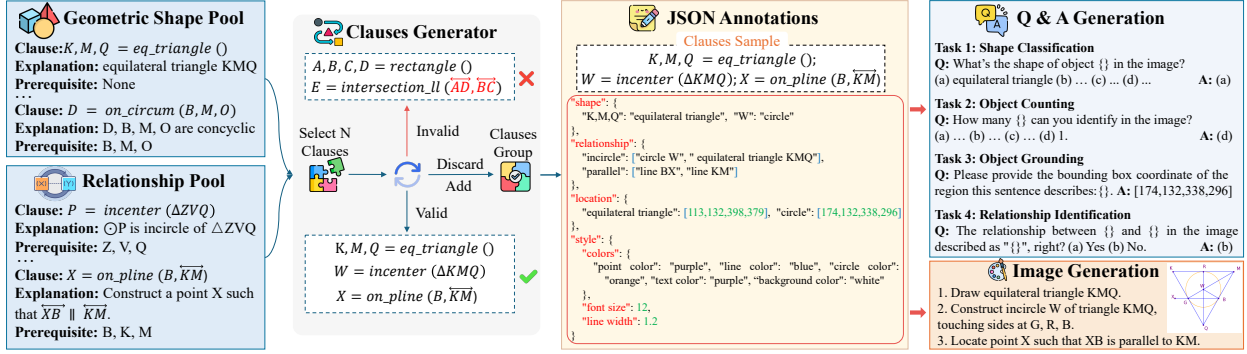
4

Figure 4. The synthetic construction process for plane geometry. We synthesize geometric figures by randomly sampling elements from the geometric shape pool and relationship pool, ensuring consistency through a verifier that enforces logical constraints based on manually designed rules, fundamental mathematical principles, and prerequisite points. All visual elements are structured and saved in JSON format. Images are rendered using the Matplotlib package, and corresponding Q&A pairs are generated using a template-based pipeline.

save the outputs as structured JSON annotations for image rendering using the Matplotlib package [22] and generating question-answer pairs via template-based pipelines.

**Image Generation.** The JSON annotations define foreground and background styles, *i.e.*, colors sampled from a monochromatic palette, line width, and font size, as well as object shape information (class names), attributes (vertices labeled with random letters), bounding box locations, and mathematical relationships. Notably, spatial relationships are generated based on the bounding box locations of two objects (*e.g.*, top-left, top-right, bottom-left, bottom-right). Similar to the image rendering process in [43], we translate geometric clauses into visual representations using Python code. This process first determines the coordinates of each point defined in the clause for basic shape visualization, and then new primitives are generated based on the specified relationships. Additionally, we apply image augmentations to increase task difficulty, *e.g.*, adding Gaussian noise, drawing irregular scribbles, and incorporating wedge-shaped symbols for congruent angles or auxiliary lines (see Figs. 9 and 11 in the Appendix).

**Question-Answer Generation.** Based on the structured annotations, we employ a template-based pipeline to generate multiple-choice and true-false questions across four perception tasks. For example, in constructing Q&A pairs for shape classification, we first randomly select an object and its associated attributes from the given figure. We then formulate a question by filling placeholders in our carefully designed templates. For true-false questions, the ground truth answer is directly derived from the JSON annotations. For multiple-choice questions, we generate plausible distractors to challenge the model's geometric perception, combining them with the correct answer into a set of four choices. Distractors are selected from geometric candidate pools, ranging from visually similar (*e.g.*, equilateral *vs.* isosceles triangles) to visually dissimilar pairs (*e.g.*, equilateral triangle *vs.* circle). Further details on the Q&A generation templates

are provided in Tab. 7 of the Appendix.

### 2.2.2. Reformat Dataset for Solid Geometry & Graphs

We leverage the public datasets CLEVR [24] and FigureQA [25] to construct benchmarks for solid geometry and graphs, respectively. Since their original annotation formats are incompatible with our template-based Q&A generation pipeline, we reformat them into structured JSON, consistent with the format used in our synthetic data process. Once reformatted, we apply the same Q&A generation strategy as used for plane geometry. See §D for details.

After generating our synthetic data and collecting public datasets, we conduct a comprehensive review to verify answer accuracy, ensure consistency between questions and diagrams, and confirm relevance to the four perception tasks, ensuring high-quality and precise dataset annotations.

### 2.3. Perception-oriented Training Dataset (GeoPeP)

Evaluation results on MATHGLANCE reveal that both open-source and closed-source MLLMs struggle to identify relevant visual regions in symbolic and abstract mathematical diagrams, despite their strong performance in theoretical reasoning and numerical computation. This contrasts with human cognitive abilities, where low-level perceptual tasks are typically solved rapidly compared to high-level reasoning tasks. To overcome such perceptual limitations, we construct a perception-oriented training dataset, GeoPeP, which enhances geometric understanding and visual grounding, ultimately benefiting complex reasoning. GeoPeP, comprising 200K high-quality image-text pairs in both diagram-caption and instruction-conversation formats, provides detailed geometric primitive information, including shapes, vertices, relationships, and fine-grained bounding box coordinates.

Unlike semantically rich natural images, mathematical diagrams are inherently abstract with structured symbols, shapes, and interrelationships. Learning from explicit struc-

Table 2. Performance comparison of different MLLMs on MATHGLANCE across plane geometry, solid geometry, and graphs. *cls*, *cnt*, *grd*, and *rlat* represent different question categories: shape classification, object counting, object grounding, and relationship identification, respectively. *all* indicates the overall accuracy, calculated as the ratio of correctly answered questions to the total number of questions in the benchmark, while **Avg.** denotes the average *all* score across all subjects.

| Model | Size | Avg. | Plane Geometry | | | | | Soild Geometry | | | | | Graphs | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | all | cls | cnt | grd | rlat | all | cls | cnt | grd | rlat | all | cls | cnt | grd | rlat |
| *Open-Source Generic MLLMs* | | | | | | | | | | | | | | | | | |
| LLaVA-v1.5 [32] | 7B | 33.3 | 29.2 | 29.0 | 39.6 | 14.2 | 37.5 | 31.6 | 43.0 | 42.3 | 0.0 | 31.3 | 39.0 | 76.8 | 35.2 | 0.0 | 39.4 |
| LLaVA-v1.5 [32] | 13B | 35.4 | 32.8 | 29.3 | 40.4 | 23.5 | 42.0 | 35.9 | 60.5 | 38.1 | 0.0 | 35.0 | 37.6 | 63.8 | 42.6 | 0.0 | 45.5 |
| mPLUG-owl3 [48] | 7B | 50.0 | 36.4 | 46.7 | 41.6 | 3.9 | 58.5 | 65.3 | 95.4 | 83.5 | 0.0 | 62.5 | 48.2 | 59.4 | 77.8 | 0.0 | 66.7 |
| InternLM-XComposer2 [16] | 7B | 55.6 | 35.8 | 49.4 | 48.8 | 0.0 | 47.0 | 62.9 | 90.7 | 86.6 | 0.0 | 53.8 | 54.6 | 60.9 | 94.4 | 0.0 | 78.8 |
| Qwen2-VL [46] | 7B | 51.4 | 37.9 | 47.6 | 41.2 | 12.8 | 53.0 | 64.1 | 93.0 | 78.4 | 14.3 | 55.0 | 52.3 | 84.1 | 88.9 | 3.2 | 18.2 |
| Qwen2-VL [46] | 72B | 59.9 | 42.4 | 51.2 | 50.8 | 17.4 | 52.0 | 71.2 | 97.7 | 84.5 | 6.4 | 77.5 | 66.1 | 76.8 | 98.2 | 16.1 | 84.9 |
| Qwen2.5-VL [6] | 7B | 59.2 | 44.0 | 56.2 | 51.3 | 18.5 | 52.0 | 68.0 | 98.8 | 88.7 | 0.0 | 65.0 | 65.7 | 89.9 | 100.0 | 3.2 | 78.8 |
| DeepSeek-VL2-Tiny [47] | 3B | 32.6 | 29.5 | 45.2 | 34.4 | 4.6 | 32.0 | 39.0 | 76.7 | 32.0 | 0.0 | 37.5 | 29.4 | 39.1 | 57.4 | 0.0 | 18.2 |
| DeepSeek-VL2-Small [47] | 16B | 51.5 | 37.6 | 47.6 | 43.6 | 12.5 | 48.5 | 63.8 | 98.8 | 70.1 | 11.1 | 60.0 | 53.2 | 76.8 | 53.7 | 11.3 | 81.8 |
| InternVL2 [11] | 8B | 48.4 | 31.9 | 44.3 | 38.0 | 0.0 | 48.5 | 62.9 | 98.8 | 62.9 | 4.8 | 70.0 | 50.5 | 68.1 | 75.9 | 0.0 | 66.7 |
| InternVL2.5 [10] | 8B | 50.7 | 35.0 | 48.8 | 36.0 | 0.0 | 60.0 | 65.6 | 98.8 | 72.2 | 4.8 | 70.0 | 51.4 | 68.1 | 77.8 | 0.0 | 69.7 |
| InternVL2.5 [10] | 38B | 63.1 | 44.0 | 59.9 | 52.0 | 2.5 | 66.0 | 78.8 | 98.8 | 92.8 | 38.1 | 72.5 | 66.5 | 98.6 | 96.3 | 3.2 | 69.7 |
| *Closed-Source Generic MLLMs* | | | | | | | | | | | | | | | | | |
| GPT-4o | - | 53.3 | 42.8 | 58.4 | 53.2 | 1.1 | 62.5 | 60.7 | 72.1 | 84.5 | 1.6 | 66.3 | 56.4 | 92.8 | 72.2 | 1.6 | 57.6 |
| GPT-o1 | - | 36.5 | 15.8 | 33.2 | 11.6 | 0.0 | 14.0 | 41.4 | 75.6 | 52.6 | 0.0 | 23.8 | 52.3 | 82.6 | 81.5 | 0.0 | 39.4 |
| *Open-Source Mathematical MLLMs* | | | | | | | | | | | | | | | | | |
| Math-LLaVA [41] | 13B | 40.0 | 27.9 | 34.4 | 32.4 | 0.0 | 50.5 | 44.8 | 81.4 | 55.7 | 0.0 | 27.5 | 47.3 | 78.3 | 59.3 | 0.0 | 51.5 |
| G-LLaVA [18] | 7B | 30.3 | 25.6 | 27.8 | 41.2 | 0.4 | 38.0 | 32.3 | 45.4 | 38.1 | 4.8 | 32.5 | 33.9 | 58.0 | 37.0 | 0.0 | 42.4 |
| MultiMath [38] | 7B | 42.1 | 31.2 | 44.0 | 30.4 | 1.07 | 53.0 | 46.7 | 81.4 | 53.6 | 4.7 | 33.8 | 48.6 | 79.7 | 57.4 | 3.2 | 33.8 |
| SVE-Math-DeepSeek [55] | 7B | 46.6 | 35.4 | 52.4 | 36.0 | 3.56 | 51.0 | 49.4 | 77.9 | 62.9 | 1.5 | 41.3 | 55.1 | 81.2 | 75.9 | 1.6 | 69.7 |
| **SVE-Math-DeepSeek$^+$ (ours)** | 7B | 68.4 | 84.6 | 75.8 | 88.4 | 82.9 | 97.5 | 54.1 | 85.3 | 65.8 | 20.3 | 45.0 | 60.7 | 85.1 | 78.4 | 1.6 | 75.7 |

tures within the data context would reduce learning complexity and enhance the problem-solving abilities of models [20]. We hypothesize that developing a structured visual dataset—where the training corpus is built from object attributes and extends to other objects based on their relationships—can improve visual attention and mitigate1 MLLMs' reliance on textual shortcuts. Thus, GeoPeP explicitly provides structured representations of fundamental geometric elements and their relationships: First, I count $\{N\}$ prominent object(s) in the image. Next, for shape information, object $\{attrbu.^i\}$ is a $\{shape^i\}$, and $\cdots$. Furthermore, I also know the fine-grained bounding box coordinates: the $\{shape^i\}$ $\{attrbu.^i\}$ is located at $\{box\_cor.^i\}$, and $\cdots$. Finally, let me explain the relationships: the $\{shape^i\}$ $\{attrbu.^i\}$ $\{rela.^{ij}\}$ to the $\{shape^j\}$ $\{attrbu.^j\}$, and $\cdots$. To further enhance the model's ability to follow instructions, we construct a task-specific instruction dataset in a multi-turn conversation format. Each question is tailored to a specific perception task, with answers presented either in free-form or as a selected option. For example: **Q**: What is the shape of object $\{attrbu.^i\}$? **A**: $\{shape^i\}$. See §E in the Appendix for additional demonstrations.

Overall, GeoPeP contributes to MLLM training by: (1) providing clear object attributes and their relationships, akin to graph nodes and edges, in training QA pairs; (2) offering fine-grained bounding box coordinates of elements, en-

abling models to systematically learn spatial awareness; and (3) integrating with reasoning-based CoT mathematical visual datasets during self-fine-tuning stage, allowing models to both perceive and reason accurately.

## 3. Experiments

In this section, we first describe the experimental setup (§3.1). We then present a comprehensive evaluation of 15 recent MLLMs (§3.2), demonstrating that while humans achieve high accuracy, MATHGLANCE remains challenging for existing models. Finally, we conduct comprehensive ablations and analyses, including detailed comparisons with math-specific MLLMs on both reasoning and perception benchmarks; investigation of key factors affecting MLLMs' perception during inference; and evaluation of how GeoPeP improves both perception and reasoning (§3.3).

### 3.1. Experimental Setup

**Generic and Mathematical MLLMs.** We evaluate 15 MLLMs on MATHGLANCE across plane geometry, solid geometry, and graphs, including closed-source generic models such as GPT-4o and GPT-o1, as well as open-source models like LLaVA-v1.5 [32], mPLUG-Owl3 [48], InternLM-XComposer2 [16], Qwen2VL [46], Qwen2.5VL [6], DeepSeek-VL2 [47], InternVL2 [11], and InternVL2.5 [10]. Additionally, we assess math-specific MLLMs, in-

cluding SVE-Math-DeepSeek [55], Math-LLaVA [41], G-LLaVA [18], and MultiMath [38]. This comprehensive evaluation provides insights into the diagram perception capabilities of state-of-the-art multimodal models.

**Implementation Details.** For open- and closed-source generic MLLMs, we follow the official inference settings, including temperature, number of beams, and maximum token length. For open-source mathematical LLMs, we adopt the standard configurations from SVE-Math-DeepSeek, setting the temperature to 0, the number of beams to 1, and the maximum token length to 1024. Notably, MultiMath [38] uses a temperature of 0.2; therefore, for a fair comparison, we rerun the inference with the temperature set to 0. We retain the original 512×512 image resolution throughout all experiments. Choices are extracted using predefined rules tailored to each MLLM's output format.

## 3.2. Main Results

The performance comparison of current MLLMs on plane geometry, solid geometry [24], and graphs [25] is summarized in Table 2. To assess models' geometric diagram perception abilities, we evaluate their performance across shape classification, object counting, relationship identification, and object grounding. Key findings reveal that both generic and mathematical MLLMs face significant challenges in geometric perception tasks.

**Generic MLLMs.** General-purpose models trained on diverse datasets, including tables, charts, and documents [9, 25, 51], as well as visual grounding datasets [40, 49], still perform poorly on mathematical diagram perception. In particular, their performance on plane geometry remains low, with most models scoring below 45% on average. For solid geometry and graphs, general-purpose models significantly outperform mathematical MLLMs, due to their exposure to large-scale FigureQA [25], CLEVR [24], and various chart understanding datasets [51]. However, they still fail in fine-grained box-level tasks, with most models achieving 0 accuracy, and lag significantly behind human-level perception. These evaluation results highlight the limitations of generic MLLMs in mathematical diagram perception. Despite being trained on large-scale and diverse natural image datasets (*e.g.*, over 2T tokens for Qwen2.5-VL and DeepSeek-VL2), these models struggle with the abstract nature of mathematical diagrams, which are defined by precise geometric structures and symbolic notations. Unlike natural images—where models can exploit semantic priors (*e.g.*, people typically appear on the ground, birds in the sky)—mathematical diagrams lack such contextual cues. Their symbolic and structured nature requires genuine visual understanding rather than superficial pattern recognition. Without this, models fail to identify where to look, leading to poor perception performance. Notably, scaling up model size is neither an optimal nor an effective

Table 3. Performance comparison of mathematical MLLMs on both math perception and reasoning benchmarks.

| Model | MATHGLANCE | | | MathVerse | MathVista | GeoQA |
|---|---|---|---|---|---|---|
| | Plane | Soild | Graphs | | | |
| Math-LLaVA [41] | 27.9 | 44.8 | 47.3 | 20.1 | 46.6 | 60.7 |
| G-LLaVA [18] | 25.6 | 31.3 | 33.9 | 17.8 | 25.6 | 64.2 |
| MultiMath [38] | 31.2 | 45.7 | 48.6 | 25.9 | 49.3 | 74.1 |
| SVE-Math-DeepSeek [55] | 35.4 | 49.4 | 55.1 | 24.3 | 48.7 | 72.8 |
| SVE-Math-DeepSeek+ | 84.6 | 54.1 | 60.7 | 28.1 | 51.3 | 76.2 |

solution, as it provides only marginal gains in perception compared to reasoning benchmarks. For instance, increasing Qwen2VL from 7B to 72B improves top-1 accuracy by 22.3% on MathVista but only 8.3% on MATHGLANCE.

**Mathematical MLLMs.** Models designed for mathematical reasoning, such as MultiMath, Math-LLaVA, and SVE-Math-DeepSeek, exhibit strong reasoning performance on benchmarks like MathVerse [53], MathVista [34], and GeoQA [18]. They struggle with fundamental geometric perception, particularly fine-grained object localization, achieving near-zero accuracy at the default IoU threshold of 0.65. Since SVE-Math-DeepSeek incorporates a geometric visual encoder trained with detection and boundary segmentation losses, it outperforms other mathematical MLLMs on our benchmark across both cross-level and fine-level perception tasks. Furthermore, as it is not exposed to solid geometry diagrams or graphical images during MLLM training stages, it provides a more controlled experimental setting for analyzing perception capabilities compared to general-purpose models. Therefore, we select SVE-Math-DeepSeek as our base model to explore enhancements in both perception and reasoning. The comparison between SVE-Math-DeepSeek and other mathematical MLLMs is shown in Tab. 3. SVE-Math-DeepSeek trained on GeoPeP significantly improves solid geometry perception. Even without direct training on solid geometry and graphs, our model outperforms others due to its enhanced ability to discriminate relationships and understand spatial configurations. Low-level diagram structure understanding is essential for high-level reasoning skills, contributing to ∼ +4% across three mathematical reasoning benchmarks. We leave the exploration of structure-aware training samples for solid geometry and graphs to future work.

## 3.3. Ablation Study

### 3.3.1. Key Factors Influencing Diagram Perception

We examine five key factors: the number of objects, object quality, visual distractors, text distractors, and Chain-of-Thought (CoT) responses. While those factors are by no means exhaustive, they aim to illuminate some of the fundamental perceptual limitations of current MLLMs, providing insights for practical applications and future improvements. The evaluation results on plane geometry from MATHGLANCE across nine MLLMs are presented
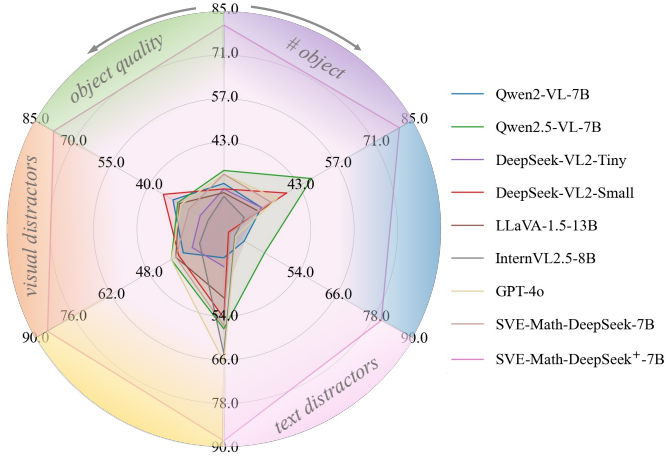
Figure 5. We evaluate five key factors affecting the perception ability of MLLMs: the number of obejcts (# obejct), visual quality, visual and textual distractors, and Chain-of-Thought (CoT) reasoning strategies (see Fig. 6 for details).
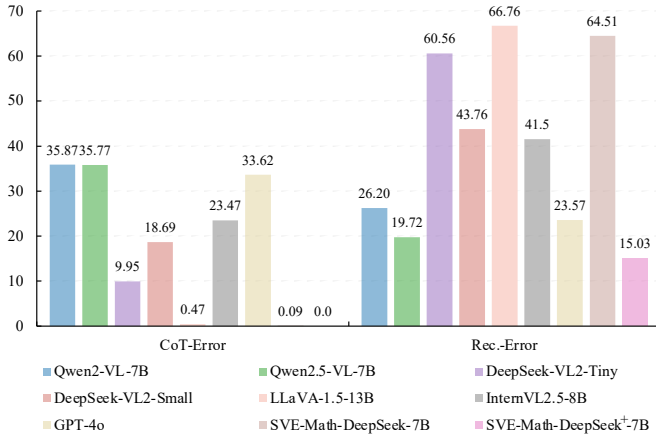


Figure 6. Analyzing perceptual errors: CoT *vs*. direct reasoning.

in Fig.5, and response error patterns are shown in Fig.6.

**The Number of Objects.** The complexity of visual perception tasks increases with the number of objects in an image. We observe that as the number of geometric shapes grows, models struggle with accurate object localization, shape classification, and counting. Fig. 1 presents the overall results across all perception tasks for eight MLLMs and our method. When scaling the number of objects from 1 to $N$ ($N < 5$), the Qwen-series model's performance drops by approximately 12%, the DeepSeek-series by 18%, and GPT-4o by 15%. Our method and its baseline show declines of 6% and 13%.

**Visual Distractors.** We draw visual distractors—irregular scribbles, wedge-shaped angle markers, and auxiliary lines—to evaluate MLLMs' ability to focus on relevant geometric elements (see Fig.9). Despite prompting models to ignore these distractors, most are still negatively affected. Close-up results are shown in Tab. 4. All evaluated MLLMs

Table 4. Accuracy (%) on plane geometry perception w.r.t. object quality (w/ or w/o Gaussian noise) and visual distractor injection.

| **MathGlance** | Noise $\sim \mathcal{N}(0, 0.3)$ | | Distractors | |
| --- | --- | --- | --- | --- |
| | w/o | w/$_{\Delta\uparrow\downarrow}$ | w/o | w/$_{\Delta\uparrow\downarrow}$ |
| Qwen2VL-7B | 29.8 | 33.9$_{\uparrow 4.1}$ | 29.8 | 26.2$_{\downarrow 3.6}$ |
| Qwen2.5VL-7B | 33.9 | 32.1$_{\downarrow 1.8}$ | 33.9 | 31.0$_{\downarrow 2.9}$ |
| DeepSeek-VL2-Tiny | 27.9 | 23.8$_{\downarrow 4.1}$ | 27.9 | 22.6$_{\downarrow 5.3}$ |
| DeepSeek-VL2-Small | 27.9 | 37.5$_{\uparrow 9.6}$ | 27.9 | 25.9$_{\downarrow 2.0}$ |
| LLaVa-1.5-13B | 26.8 | 31.6$_{\uparrow 4.8}$ | 26.8 | 29.1$_{\uparrow 2.3}$ |
| InternVL2.5-8B | 25.6 | 20.8$_{\downarrow 4.8}$ | 25.6 | 19.6$_{\downarrow 6.0}$ |
| GPT-4o | 32.7 | 30.9$_{\downarrow 1.8}$ | 32.7 | 28.6$_{\downarrow 4.1}$ |
| SVE-Math-DeepSeek-7B | 32.8 | 28.0$_{\downarrow 4.8}$ | 32.8 | 29.2$_{\downarrow 3.6}$ |
| SVE-Math-DeepSeek$^+$-7B | 80.7 | 78.3$_{\downarrow 2.4}$ | 80.7 | 80.3$_{\downarrow 0.4}$ |

Table 5. Performance of relationship identification (*rlat*) on MathGlance under different textual distractor settings.

| **MathGlance** (*rlat*) | Unrela. infor. | | Conflicts | |
| --- | --- | --- | --- | --- |
| | w/o | w/$_{\Delta\uparrow\downarrow}$ | w/o | w/$_{\Delta\uparrow\downarrow}$ |
| Qwen2VL-7B | 53.0 | 48.5$_{\downarrow 4.5}$ | 53.0 | 24.5$_{\downarrow 28.5}$ |
| Qwen2.5VL-7B | 52.0 | 56.0$_{\uparrow 4.0}$ | 52.0 | 30.0$_{\downarrow 22.0}$ |
| DeepSeek-VL2-Tiny | 32.0 | 34.5$_{\uparrow 2.5}$ | 32.0 | 3.5$_{\downarrow 28.5}$ |
| DeepSeek-VL2-Small | 48.5 | 42.5$_{\downarrow 6.0}$ | 48.5 | 20.5$_{\downarrow 28.0}$ |
| LLaVa-1.5-13B | 42.0 | 37.5$_{\downarrow 4.5}$ | 42.0 | 16.0$_{\downarrow 26.0}$ |
| InternVL2.5-8B | 60.0 | 49.0$_{\downarrow 11.0}$ | 60.0 | 18.0$_{\downarrow 42.0}$ |
| GPT-4o | 62.5 | 63.5$_{\uparrow 1.0}$ | 62.5 | 38.0$_{\downarrow 24.5}$ |
| SVE-Math-DeepSeek-7B | 51.0 | 49.0$_{\downarrow 2.0}$ | 51.0 | 15.0$_{\downarrow 36.0}$ |
| SVE-Math-DeepSeek$^+$-7B | 88.0 | 86.5$_{\downarrow 1.5}$ | 88.0 | 74.0$_{\downarrow 14.0}$ |

show a performance drop of 2–10% under visual distractors, except LLaVA-1.5-13B, which improves by 2.3%. Our model remains robust, showing minimal sensitivity to distractors and better visual focus on geometric primitives.

**Object Quality.** For visual fidelity analysis, we apply Gaussian noise to degrade object quality. As shown in Tab. 4, with a default variance of 0.3, most models show performance degradation, while Qwen2VL, DeepSeek-VL2-Small, and LLaVA-1.5-13B unexpectedly improve. This may stem from their exposure to visually degraded data during training, making noisy images more aligned with their learned distribution—especially in the absence of clean geometric diagrams in the training corpus. To enable a more comprehensive evaluation and mitigate data bias, we introduce large Gaussian noise with variances of 0.5/0.8 to assess model performance under severe visual distortion. As noise increases, all models struggle with geometric recognition, showing significant performance drops. Fig.10 reports perception accuracy under noise, and distorted examples are shown in Appendix Fig.11, where at variance 0.8, even humans fail to interpret the diagrams.

**Text Distractors.** Researchers have observed that large language models (LLMs) often dominate MLLMs, relying heavily on knowledge shortcuts rather than genuinely

**Question:**
Two triangles are similar. WS is 8, and VT is 10. Find WR.

**Choices:**
A: 6    **B: 8**
C: 10    D: 12

**SVE-Math-DeepSeek⁺**

✓ Since triangle WSR is similar to triangle VTR, we have $\frac{WS}{VT} = \frac{RW}{TR}$.

✓ Given that WS = 8 and VT = 10, we can calculate RW as $RW = \frac{WS \times TR}{VT} = \frac{8 \times 10}{10} = 8$. So, the correct answer is **B**.

**SVE-Math-DeepSeek**

✗ Using the property of similar triangles, we can establish the equation $\frac{WS}{VT} = \frac{WR}{VS}$.

✗ Substituting the given values, we get $\frac{8}{10} = \frac{WR}{6}$. Solving for WR yields $WR = \frac{48}{10}$. Therefore, the answer is option **D**.

**InternVL2.5**

✗ Since the triangles are similar, the ratio of the corresponding sides is equal, so $\frac{WS}{VT} = \frac{WR}{VR}$.

✗ Substitute the given values: $\frac{8}{10} = \frac{WR}{WR+6}$. Then, WR=24. Therefore, the correct answer is **A**.

Figure 7. Response comparisons between SVE-Math-DeepSeek⁺, SVE-Math-DeepSeek, and InterVL2.5.

perceiving visual content. To investigate this, we evaluate two settings in the relationship identification task. The first involves introducing unrelated information w.r.t. the target task. For example, when the task is to identify a mathematical relationship, the model is additionally given irrelevant details about foreground/background colors, box shapes, or object attributes. The second setting introduces contradictory cues; for instance, if two lines are parallel, the model is explicitly prompted that they are perpendicular. The results in Tab. 5 show that providing conflicting knowledge significantly impairs the perceptual ability of all evaluated MLLMs. The impact is especially pronounced in models with weaker perception capabilities, such as InternVL2.5-8B (48.4% in avg.), while stronger perceptual models like Qwen2.5-VL-7B (59.2% in avg.) are less affected. Our method demonstrates greater robustness, though it still experiences a 14% performance drop. These findings highlight that reliance on knowledge shortcuts remains a critical issue, particularly in abstract diagram recognition, where models must genuinely understand visual content rather than rely on superficial pattern matching.

**Chain-of-Thought (CoT) Response.** CoT reasoning is designed to enhance step-by-step logical inference, yet its effectiveness in visual perception tasks remains uncertain. Our analysis shows that while CoT improves textual reasoning, it does not directly enhance spatial or geometric understanding. Models incorporating CoT reasoning often struggle with fundamental perception tasks, leading to significant CoT errors. Models often generate excessive yet irrelevant rationale misaligned with the diagram, ultimately resulting in incorrect responses, particularly in Qwen-series models and GPT-4o (over 30% in Fig. 6). See §G for CoT response examples in the Appendix.

### 3.3.2. Effect of GeoPeP

Building on SVE-Math-DeepSeek, we first train the projector while freezing the LLM and visual encoder during the visual-language alignment stage using either AutoGeo, MAVIS, or GeoPeP image-caption alignment datasets. As shown in Tab. 6, training with GeoPeP yields a 6.2% gain on perception tasks, but suffer a substantial drop on the other two datasets (~-27%). Using these three pre-

Table 6. Performance comparison w.r.t. variants trained on GeoPeP and other mathematical visual alignment datasets. ★ indicates the pretrained model continuously self-fine-tuned using MathV360K and instruction-formatted GeoPeP.

| Model | | MATHGLANCE | | | MathVerse | MathVista |
|---|---|---|---|---|---|---|
| | | Plane | Soild | Graphs | | |
| SVE-Math-DeepSeek | | 35.4 | 49.4 | 55.1 | 24.3 | 48.7 |
| Align. | +AutoGeo (△) | 18.6 | 34.7 | 40.4 | 17.8 | 45.0 |
| | +MAVIS (◇) | 6.0 | 2.5 | 6.9 | 6.8 | 34.3 |
| | +GeoPeP (♡) | 41.6 | 40.2 | 49.2 | 19.7 | 46.2 |
| SFT | △ ● ★ | 79.6 | 44.2 | 55.1 | 25.2 | 48.3 |
| | ◇ ● ★ | 78.1 | 42.3 | 54.1 | 23.3 | 47.1 |
| | ♡ ● ★ | **84.6** | **54.1** | **60.7** | **28.1** | **51.3** |

trained models, we then conduct self-fine-tuning (SFT) with MathV360K and instruction-formatted GeoPeP. The model trained with GeoPeP achieves a notable ~4% improvement on both MathVerse and MathVista compared to SVE-Math-DeepSeek, whereas other variants underperform on Math-Vista. We hope this observation encourages further research into more effective methods and training datasets to enhance the visual perception capabilities of mathematical MLLMs, as clearer visual understanding enables more faithful and accurate reasoning responses (Figs.7/15-16).

## 4. Conclusion

In this paper, we introduce MATHGLANCE, a perception-demanding benchmark designed to evaluate geometric perception in MLLMs. While these tasks appear trivial for humans to solve 'at a glance', they present significant challenges for current models. Unlike semantic-rich natural images, mathematical diagrams are structured and abstract, characterized by symbolic elements and their precise interrelationships. Thus, they require genuine visual interpretation rather than pattern recognition or textual shortcuts. Through extensive analysis, we identify key factors affecting diagam perception, highlighting the importance of high-quality and structure-aware training datasets to enable MLLMs to both see and reason effectively. Our findings suggest that improving low-level geometric perception is essential for advancing high-level reasoning, bridging the gap between visual understanding and logical inference.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1

[2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1

[3] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *Advances in neural information processing systems*, pages 23716–23736, 2022. 1

[4] Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. Mathqa: Towards interpretable math word problem solving with operation-based formalisms. In *Annual Meeting of the Association for Computational Linguistics*, pages 2357–2367, 2019. 1

[5] Avinash Anand, Raj Jaiswal, Abhishek Dharmadhikari, Atharva Marathe, Harsh Popat, Harshil Mital, Ashwin R Nair, Kritarth Prasad, Sidharth Kumar, Astha Verma, et al. Geovqa: A comprehensive multimodal geometry dataset for secondary education. In *2024 IEEE 7th International Conference on Multimedia Information Processing and Retrieval (MIPR)*, pages 102–108. IEEE, 2024. 1

[6] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 2, 6

[7] Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou. Hallucination of multimodal large language models: A survey. *arXiv preprint arXiv:2404.18930*, 2024. 2

[8] Jiaqi Chen, Tong Li, Jinghui Qin, Pan Lu, Liang Lin, Chongyu Chen, and Xiaodan Liang. Unigeo: Unifying geometry logical reasoning via reformulating mathematical expression. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 3313–3323, 2022. 1

[9] Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyou Zhou, and William Yang Wang. Tabfact: A large-scale dataset for table-based fact verification. *arXiv preprint arXiv:1909.02164*, 2019. 7

[10] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024. 6, 3

[11] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei 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. 6

[12] An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision-language models. *Advances in Neural Information Processing Systems*, 37:135062–135093, 2025. 1

[13] J. G. Cromley, L. E. Snyder-Hogan, and U. A. Luciw-Dubas. Cognitive activities in complex science text and diagrams. *Contemporary Educational Psychology*, 35(1):59–74, 2010. 2

[14] M. de Rijke. Logical reasoning with diagrams. *Journal of Logic, Language and Information*, 8(3):387–390, 1999. 2

[15] Linger Deng, Yuliang Liu, Bohan Li, Dongliang Luo, Liang Wu, Chengquan Zhang, Pengyuan Lyu, Ziyang Zhang, Gang Zhang, Errui Ding, et al. R-cot: Reverse chain-of-thought problem generation for geometric reasoning in large multimodal models. *arXiv preprint arXiv:2410.17885*, 2024. 1

[16] Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang Zhang, Haodong Duan, Maosong Cao, et al. Internlm-xcomposer2: Mastering free-form text-image composition and comprehension in vision-language large model. *arXiv preprint arXiv:2401.16420*, 2024. 6

[17] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. In *International Conference on Machine Learning*, pages 8469–8488, 2023. 1

[18] Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, Lanqing Hong, Jianhua Han, Hang Xu, Zhenguo Li, et al. G-llava: Solving geometric problem with multi-modal large language model. *arXiv preprint arXiv:2312.11370*, 2023. 6, 7, 1

[19] Yash Goyal, Tejas Khot, Aishwarya Agrawal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. *International Journal of Computer Vision*, 127:398–414, 2019. 1

[20] Haoyu Han, Yaochen Xie, Hui Liu, Xianfeng Tang, Sreyashi Nag, William Headden, Yang Li, Chen Luo, Shuiwang Ji, Qi He, et al. Reasoning with graphs: Structuring implicit knowledge to enhance llms reasoning. *arXiv preprint arXiv:2501.07845*, 2025. 6

[21] Zihan Huang, Tao Wu, Wang Lin, Shengyu Zhang, Jingyuan Chen, and Fei Wu. Autogeo: Automating geometric image dataset creation for enhanced geometry understanding. *arXiv preprint arXiv:2409.09039*, 2024. 2, 1

[22] John D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science Engineering*, 9(3):90–95, 2007. 5

[23] Chaoya Jiang, Haiyang Xu, Mengfan Dong, Jiaxing Chen, Wei Ye, Ming Yan, Qinghao Ye, Ji Zhang, Fei Huang, and Shikun Zhang. Hallucination augmented contrastive learning for multimodal large language model. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 27036–27046, 2024. 2

[24] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2901–2910, 2017. 3, 4, 5, 7, 2

[25] Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. Figureqa: An annotated figure dataset for visual reasoning. *arXiv preprint arXiv:1710.07300*, 2017. 3, 4, 5, 7, 2

[26] Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. Seed-bench: Benchmarking multimodal large language models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 13299–13308, 2024. 1

[27] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In *The Conference on Empirical Methods in Natural Language Processing*, 2023. 1

[28] Zechao Li, Yanpeng Sun, Liyan Zhang, and Jinhui Tang. Ct-net: Context-based tandem network for semantic segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(12):9904–9917, 2021. 1

[29] Zhenwen Liang, Tianyu Yang, Jipeng Zhang, and Xiangliang Zhang. Unimath: A foundational and multimodal mathematical reasoner. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 7126–7133, 2023. 1

[30] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755, 2014. 1

[31] Fangyu Liu, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Yasemin Altun, Nigel Collier, and Julian Martin Eisenschlos. Matcha: Enhancing visual language pretraining with math reasoning and chart derendering. In *The Annual Meeting Of The Association For Computational Linguistics*, 2023. 1

[32] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances in neural information processing systems*, pages 34892–34916, 2023. 1, 6

[33] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European Conference on Computer Vision*, pages 216–233, 2024. 1

[34] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *International Conference on Learning Representations*. 2, 7, 1

[35] Pan Lu, Ran Gong, Shibiao Jiang, Liang Qiu, Siyuan Huang, Xiaodan Liang, and Song-chun Zhu. Inter-gps: Interpretable geometry problem solving with formal language and symbolic reasoning. In *Annual Meeting of the Association for Computational Linguistics*, pages 6774–6786, 2021. 1

[36] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong Wang. Conditional detr for fast training convergence. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3651–3660, 2021. 1

[37] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual question answering by reading text in images. In *International Conference on Document Analysis and Recognition*, pages 947–952. IEEE, 2019. 1

[38] Shuai Peng, Di Fu, Liangcai Gao, Xiuqin Zhong, Hongguang Fu, and Zhi Tang. Multimath: Bridging visual and mathematical reasoning for large language models. *arXiv preprint arXiv:2409.00147*, 2024. 3, 6, 7

[39] Minjoon Seo, Hannaneh Hajishirzi, Ali Farhadi, Oren Etzioni, and Clint Malcolm. Solving geometry problems: Combining text and diagram interpretation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1466–1476, 2015. 1

[40] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 8429–8438, 2019. 7

[41] Wenhao Shi, Zhiqiang Hu, Yi Bin, Junhua Liu, Yang Yang, See-Kiong Ng, Lidong Bing, and Roy Ka-Wei Lee. Math-llava: Bootstrapping mathematical reasoning for multimodal large language models. *arXiv preprint arXiv:2406.17294*, 2024. 3, 6, 7

[42] Yanpeng Sun, Huaxin Zhang, Qiang Chen, Xinyu Zhang, Nong Sang, Gang Zhang, Jingdong Wang, and Zechao Li. Improving multi-modal large language model through boosting vision capabilities. *arXiv preprint arXiv:2410.13733*, 2024. 1

[43] Trieu Trinh, Yuhuai Tony Wu, Quoc Le, He He, and Thang Luong. Solving olympiad geometry without human demonstrations. *Nature*, 625:476–482, 2024. 4, 5, 1

[44] Chenxi Wang, Xiang Chen, Ningyu Zhang, Bozhong Tian, Haoming Xu, Shumin Deng, and Huajun Chen. Mllm can see? dynamic correction decoding for hallucination mitigation. In *The Thirteenth International Conference on Learning Representations*, 2024. 2

[45] Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. In *Advances in Neural Information Processing Systems*, pages 95095–95169, 2025. 1

[46] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 6, 1

[47] Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, et al. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding. *arXiv preprint arXiv:2412.10302*, 2024. 6, 1, 2

11

[48] Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl3: Towards long image-sequence understanding in multimodal large language models. In *International Conference on Learning Representations*, 2024. 6

[49] Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity. In *The Twelfth International Conference on Learning Representations*, 2024. 7

[50] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. In *International Conference on Machine Learning*, pages 57730–57754, 2024. 1

[51] Ye Yuan, Xiao Liu, Wondimu Dikubab, Hui Liu, Zhilong Ji, Zhongqin Wu, and Xiang Bai. Syntax-aware network for handwritten mathematical expression recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4553–4562, 2022. 7

[52] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024. 1

[53] Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu 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. 2, 7, 1

[54] Renrui Zhang, Xinyu Wei, Dongzhi Jiang, Ziyu Guo, Shicheng Li, Yichi Zhang, Chengzhuo Tong, Jiaming Liu, Aojun Zhou, Bin Wei, et al. Mavis: Mathematical visual instruction tuning with an automatic data engine. *arXiv preprint arXiv:2407.08739*, 2024. 2, 1

[55] Shan Zhang, Aotian Chen, Yanpeng Sun, Jindong Gu, Yi-Yu Zheng, Piotr Koniusz, Kai Zou, Anton van den Hengel, and Yuan Xue. Open eyes, then reason: Fine-grained visual mathematical understanding in mllms. *arXiv preprint arXiv:2501.06430*, 2025. 2, 6, 7, 3

[56] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations*, 2020. 1

# MATHGLANCE: Multimodal Large Language Models Do Not Know Where to Look in Mathematical Diagrams

## Appendix

## A. Related Work

### A.1. Mathematical Reasoning Benchmark

To evaluate MLLMs performance across different domains, various benchmarks [19, 26, 27, 33, 50] have been proposed, primarily focusing on natural scene understanding. However, benchmarks specifically designed for multimodal mathematical reasoning remain scarce.

Early benchmarks such as MathQA [4], UniGeo [8]Geometry3k [35], GEOS [39], and GeoQA++ [5] introduced multimodal mathematical tasks but were limited in scope, often focusing on specific subdomains like plane geometry. More recent efforts have sought to provide broader and more diverse evaluations. MMMU [52] assesses multimodal mathematical understanding with an emphasis on symbolic reasoning and word problem-solving. MathVista [34] targets geometry-related tasks by integrating both real-world and synthetic diagrams to evaluate visual reasoning. MathVerse [53] expands this by incorporating a wider range of multimodal challenges involving charts, graphs, and structured visual content. MATH-V [45] further addresses the limitations of existing benchmarks by curating 3,040 high-quality math problems sourced from real-world competitions. While these benchmarks introduce visual elements, their core focus remains on assessing mathematical reasoning. It remains unclear whether the performance on these tasks truly reflects the models' ability to comprehensively understand the symbolic information in diagrams. The ability to accurately interpret mathematical symbols, diagrams, and spatial structures is a fundamental component of solving multimodal math problems, yet current benchmarks place limited emphasis on this aspect.

### A.2. MLLMs for Math

While multimodal large language models [3, 32, 42] have made significant strides in vision-language understanding, their mathematical reasoning capabilities remain limited. Foundation models such as GPT-4V [2], Qwen2-VL [46], and Deepseek-VL2 [47] perform well on general multimodal tasks but struggle with mathematical symbol recognition, spatial reasoning, and logical deduction, making them inadequate for vision-based math problem-solving.

To address these limitations, recent efforts have introduced math-specific MLLMs for improving mathematical reasoning. AlphaGeometry [43] achieves state-of-the-art results in geometry by leveraging theorem-proving and reinforcement learning. However, it relies solely on text-based diagram descriptions, lacking direct image processing. G-LLaVA [18] extends LLaVA with geometric reasoning capabilities but struggles with complex visual structures and generalization beyond plane geometry. UniMath [29] integrates structured math representations for solving visual word problems but remains focused on symbolic reasoning, limiting its handling of free-form mathematical diagrams. MatCha [31] specializes in chart-based reasoning, extracting quantitative relationships from structured visual data. However, its reliance on predefined formats limits adaptability to unstructured mathematical visuals.

Beyond architecture improvements, MAVIS [54] introduces an automated data engine to generate large-scale mathematical visual datasets, reducing annotation costs while ensuring high-quality diagram-caption pairs and problem-solving rationales. However, its diagrams are constructed by simply combining basic geometric shapes, lacking mathematical relationship constraints such as perpendicularity and parallelism. AutoGeo [21] considers special properties of lines, such as midlines and radii, as foundational to many geometric theorems, incorporating these properties into geometric figures. Reverse Chain-of-Thought (R-CoT) [15] introduces the Geometry Generation Chain to generate the geometry image and corresponding description. However, current data engines still fail to explore the underlying structures in mathematical diagrams, leading to ambiguous captions and redundant information that negatively impact MLLMs' perception abilities. Consequently, models trained on such datasets perform poorly on perception-demanding tasks, even falling below random guessing on multiple-choice questions. Thus, building a mathematical visual dataset that enables models to perceive and reason is a urgent task.

## B. Visualization of MATHGLANCE

MATHGLANCE is specifically designed to evaluate perception-demanding tasks in Multimodal Large Language Models (MLLMs), focusing on four core visual understanding capabilities: shape classification, object counting, object grounding,

and relationship identification. These tasks, commonly studied in classical computer vision, are typically solvable by humans with minimal cognitive effort. Our benchmark spans three key domains—plane geometry, solid geometry, and mathematical graphs—to ensure broad coverage of visual representations encountered in educational and scientific contexts. Fig. 8 provides illustrative examples from MATHGLANCE, demonstrating the diagrammatic input and corresponding question–answer (Q&A) pairs used for evaluation.

## C. Problem Templates

This section introduces the problem templates employed in MATHGLANCE, with illustrative examples provided in Tab. 7. We adopt a template-based generation engine to systematically construct diverse perception-oriented Q&A pairs. Each question is generated by parsing the structured JSON annotations of a given image, which include information on shapes, object attributes, spatial positions, and inter-object relationships. These elements are then filled into carefully designed templates (Tab. 7) to produce grammatically correct and semantically meaningful Q&A pairs.

Our template designs are tailored to accommodate diverse subject domains—including plane geometry, solid geometry, and graph-based diagrams—and span a spectrum of perceptual complexity, from coarse-level to fine-grained tasks. For example, a coarse-level shape classification template may pose a question such as, "What is the shape of the object with {vertices}?", where the correct answer is directly retrieved from the structured annotations. For fine-grained object localization, we adopt a grounding template similar to that used in [47], prompting models with: "Please provide the bounding box coordinates of the region this sentence describes: {shape} {vertices}". This allows us to evaluate the model's ability to extract precise spatial information and align language with visual primitives at a granular level.

We design three types of questions: multiple-choice, true/false, and open-ended. Specifclly, for multiple-choice questions, we incorporate plausible distractors to challenge the model's geometric perception. These distractors are selected from geometric candidate pools, varying from visually similar to clearly distinguishable options, ensuring a balanced evaluation of fine-grained visual discrimination. By leveraging this template engine, MATHGLANCE ensures consistent question formulation, controlled variation in difficulty, and robust coverage of geometric primitives and their relationships, forming a reliable testbed for evaluating diagram perception in MLLMs.

## D. Reformat Dataset for Solid Geometry & Graphs

**Solid Geometry.** CLEVR [24] is a synthetic Visual Question Answering (VQA) dataset containing 3D-rendered objects. Each object in the scene is defined by its 3D position and a set of attributes, including size (small or large), shape (cube, cylinder, sphere), material (rubber or metal), and color (gray, blue, brown, yellow, red, green, purple, cyan). However, CLEVR's annotation format is incompatible with our template-based pipeline for generating perception evaluation QA pairs. To resolve this, we reformat the original annotations into structured JSON, following the same format as the synthetic data process. Specifically, we extract object information from CLEVR's scene annotations and organize them into shape information and object attributes. We then leverage an object's 3D location and its direction for calculating its 2D top-left and bottom-right coordinates, from which relative spatial relationships are derived.

Once the reformatted JSON annotations are generated, we apply the same template-based QA generation process as designed for plane geometry, with an additional uniqueness check for shape classification. To ensure clear and consistent labeling, we filter out ambiguous cases where identical attribute combinations correspond to different shapes within the same image. For example, if an image contains multiple distinct shapes with the same attributes (*e.g.*, a large rubber blue square and a large rubber blue cylinder), we exclude classification questions that rely on the attribute combination 'large rubber blue' to prevent ambiguity.

**Graphs.** FigureQA [25] is a visual reasoning dataset containing over one million question-answer pairs, grounded in synthetic, scientific-style figures, including line plots, dot-line plots, vertical and horizontal bar graphs, and pie charts. The official annotations capture various relationships between plot elements and evaluate characteristics such as maximum, minimum, smoothness, and intersection, all framed as binary yes/no questions. Additionally, FigureQA provides numerical data annotations used to generate each figure, along with bounding-box annotations for all plot elements. Similar to CLEVR, we reformat FigureQA annotations compatible with template-based pipelines.

Specifically, the shapes include five graph types as defined in the official specification, with each element's attribute represented by its unique color. Additionally, we store the bounding box coordinates of each foreground element (such as lines, bars, and pie slices), along with legends and titles, as ground truth for the grounding task. For relationships, we follow the binary true/false question-answering problems to evaluate the model's understanding of geometric and graphical relationships.

## E. Example Illustration of GeoPeP

Figs. 12-14 demonstrate how GeoPeP delivers explicit geometric information—covering object count, shape classification, fine-grained bounding box coordinates, and inter-object relationships—presented in both caption-style and instruction-following conversational formats.

## F. Visual Distractors

To evaluate the robustness of MLLMs' perceptual capabilities, we introduce a set of visual perturbations through data augmentation. These include Gaussian noise, irregular scribbles, wedge-shaped symbols, and auxiliary lines, as shown in Fig. 9. The aim is to simulate real-world visual ambiguities and assess whether MLLMs can retain accurate geometric understanding under degraded conditions. These controlled distortions provide a rigorous benchmark for evaluating the models' ability to extract and interpret mathematical structures from visually complex inputs.

By gradually increasing the Gaussian noise level, we can systematically evaluate its impact on the ability of MLLM to recognize mathematical structures. As shown in Fig. 11, the image transitions from a noise level of 0.1 to 0.8, progressively blurring the geometric features. We observe that as the noise intensity increases, the difficulty of recognizing mathematical structures also rises, posing greater challenges to accurate recognition and reasoning. When the noise level reaches $\sim \mathcal{N}(\mathbf{0}, \mathbf{0.8})$, the mathematical structures become nearly indistinguishable to the human eye, resulting in a significant performance drop on MATHGLANCE, as illustrated in Fig. 10.

## G. Case Study

**Model Responses.** In Fig. 15 and Fig. 16, we present a comparative analysis of the performance of different Multimodal Large Language Models (MLLMs), including SVE-Math-DeepSeek [55] (baseline), SVE-Math-DeepSeek$^+$, and InterVL2.5 [10]. Our evaluation reveals that SVE-Math-DeepSeek$^+$ achieves a significant improvement in response accuracy compared to its baseline, SVE-Math-DeepSeek. This improvement suggests that our carefully curated training data effectively enhances the model's mathematical perception, enabling it to better comprehend problem structures, reason through mathematical concepts, and generate more precise answers. Furthermore, we observe that even InterVL2.5, a strong competitor in the field, produces incorrect responses in certain cases. Upon closer examination, these errors often stem from inaccurate mathematical perception rather than purely computational mistakes. This observation reinforces the notion that an MLLM's ability to accurately perceive and interpret mathematical content is a critical factor in achieving high performance. Overall, our findings highlight the fundamental role of precise perception in mathematical reasoning, akin to its importance in vision and language understanding. Just as accurate perception is essential for tasks such as image recognition and natural language processing, it is equally vital for MATH-related problem-solving. Our study underscores that enhancing a model's mathematical perception can lead to substantial gains in accuracy and reliability.

**Error Examples.** In this section, we provide more detailed error examples of GPT-4o, Qwen2.5-VL-7B, and our 7B model. We categorize the errors into two types: Chain-of-Thought (CoT) Errors and Recognition Errors. CoT errors occur when the model engages in step-by-step reasoning for perception questions but ultimately provides incorrect answers. Recognition errors, on the other hand, arise when the model attempts direct answering without reasoning yet fails to produce the correct result. Representative examples for each error type are illustrated in Figs. 17-26.

Figure 8. Visualization of sample cases from MATHGLANCE. (a), (b), and (c) correspond to problems related to Plane Geometry, Solid Geometry, and Graphs, respectively.

4

Figure 9. Visualizing distractors w.r.t. adding Gaussian noise, drawing irregular scribbles and incorporating wedge-shaped symbols or auxiliary lines.

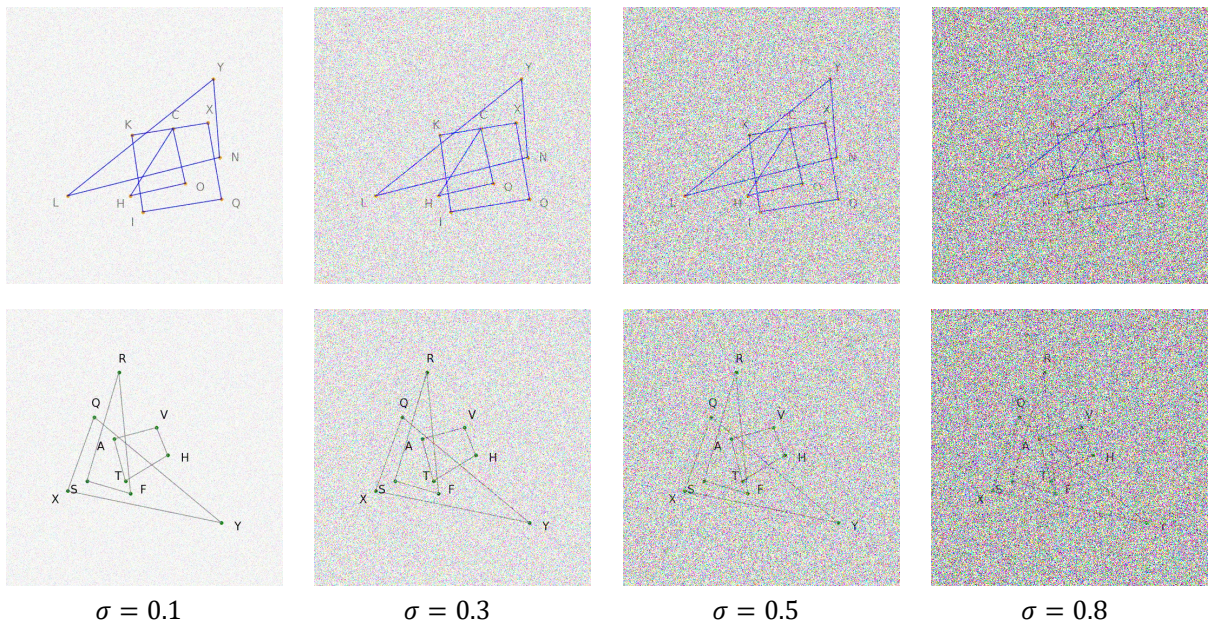Figure 10. Perceptual perfromance w.r.t. varying Gaussian noise variance.



$\sigma = 0.1$      $\sigma = 0.3$      $\sigma = 0.5$      $\sigma = 0.8$

Figure 11. Visualization of distorted images under varying levels of Gaussian noise.

## Structured Descriptions

**Image:**



**Question:** Please provide structured geometric information for the image.

**Answer:** First, I count 3 prominent objects in the image, including 1 rectangle, 1 isosceles triangle, 1 square. Next, for the shape information, object DIKP is a rectangle, object EHV is an isosceles triangle, object JNOY is a square. Furthermore, I also know the fine-grained bounding box coordinates: the rectangle DIKP is located at [113, 205, 251, 328], the isosceles triangle EHV is located at [236, 215, 399, 317], the square JNOY is located at [154, 179, 230, 256]. Finally, let me explain the relationships: rectangle DIKP is in the bottom left relative to isosceles triangle EHV; rectangle DIKP is in the bottom left relative to square JNOY.

## Multi-turn Conversation

**Image:**



**Question:** Provide the shape type of object DKLT in this image.

**Answer:** rectangle

**Question:** Count the number of rectangles in the image.

**Answer:** 2

**Question:** Please provide the bounding box coordinate of the region this sentence describes: rectangle EHIX.

**Answer:** [261, 143, 364, 251]

**Question:** Describe the relationship of rectangle EHIX relative to trapezoid FRQO in the image.

**Answer:** bottom right

Figure 12. Sampels of caption-style and instruction-following GeoPeP.

## Structured Descriptions

**Image:**



**Question:** Please provide structured geometric information for the image.

**Answer:** First, I count 1 prominent object in the image, including 1 scalene triangle. Next, for the shape information, object AHN is a scalene triangle. Furthermore, I also know the fine-grained bounding box coordinates: the scalene triangle AHN is located at [113, 164, 399, 348]. Finally, let me explain the relationships: point B lies on line NA; The angles angle GNA, angle GNE and angle ENG form a trisection of angle ANH, dividing it into three equal parts.

## Multi-turn Conversation

**Image:**



**Question:** Can you identify the type of object CGHP in the picture? A:equilateral quadrilateral B:trapezoid C:right trapezoid D:right triangle
**Answer:** B

**Question:** How many trapezoids are there in the image? A:3 B:4 C:1 D: 2
**Answer:** C

**Question:** Please provide the bounding box coordinate of the region this sentence describes: trapezoid CGHP.
**Answer:** [113, 133, 399, 380]

**Question:** Does the image show the relationship among angle CHG, angle CHD, angle DHM and angle MHG as "angle trisection"? A:Yes B:No.
**Answer:** A

Figure 13. Sampels of caption-style and instruction-following GeoPeP.

## Structured Descriptions

**Image:**



**Question:** Please provide structured geometric information for the image.

**Answer:** First, I count 2 prominent objects in the image, including 1 isosceles triangle, 1 circle. Next, for the shape information, object JNR is an isosceles triangle, object HJNR is a circle. Furthermore, I also know the fine-grained bounding box coordinates: the isosceles triangle JNR is located at [179, 143, 318, 343], the circle HJNR is located at [169, 143, 399, 369]. Finally, let me explain the relationships: isosceles triangle JNR is in the top left relative to circle HJNR; point F is the reflection of point H with respect to point R.

## Multi-turn Conversation

**Image:**



**Question:** Describe the geometric shape of object HKUOQ in the picture.
**Answer:** pentagon

**Question:** Identify the total number of pentagons present in this image.
**Answer:** 1

**Question:** Please provide the bounding box coordinate of the region this sentence describes: pentagon HKUOQ.
**Answer:** [112, 128, 399, 384]

**Question:** Describe the geometric relationship between line HE and line KO in the image.
**Answer:** perpendicular

Figure 14. Sampels of caption-style and instruction-following GeoPeP.

Table 7. Examples of problem templates used by MATHGLANCE on different source data acorss different tasks.

| Source | Task | Three randomly chosen examples from hundreds. |
|---|---|---|
| Plane Geometry | cls | What is the shape of object {vertices} in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Can you identify the type of object {vertices} in the picture? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Can you identify the type of object in the picture? Choices: A:{a} B:{b} C:{c} D:{d}.<br>What is the shape of the object in the image? Choices: A:{a} B:{b} C:{c} D:{d}. |
| | cnt | There {be} {num} {shape} in the picture, right? Choices: A:Yes B:No.<br>{be} there {num} {shape} in the picture? Choices: A:Yes B:No.<br>You can see {num} objects in the picture, can't you? Choices: A:Yes B:No.<br>There should be {num} shapes in the picture, correct? Choices: A:Yes B:No.<br>How many {shape}s can you find in the picture? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Please count all the {shape}s in the image. Choices: A:{a} B:{b} C:{c} D:{d}.<br>What is the total number of shapes in the picture? Choices: A:{a} B:{b} C:{c} D:{d}.<br>How many objects can you identify in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Please identify and select all the {shape}s in the picture. Choices: A:{a} B:{b} C:{c} D:{d}.<br>Find and select all the {shape}s in the picture. Choices: A:{a} B:{b} C:{c} D:{d}. |
| | grd | Please provide the bounding box coordinate of the region this sentence describes: {shape} {vertices}. |
| | rlat | Can the relationship {preposition} {shape} in the image be described as "{relation}"? Choices: A:Yes B:No.<br>Does the image show the relationship {preposition} {shape} as "{relation}"? Choices: A:Yes B:No.<br>Is {shape1} described as being in the '{relation}' relative to {shape2} in the image? Choices: A:Yes B:No.<br>Is {shape1} said to be in the '{relation}' relative to {shape2} in the image? Choices: A:Yes B:No.<br>What is the relationship {preposition} {shape} in the picture? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Can you identify the relationship {preposition} {shape} in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>What is the relative position of {shape1} to {shape2} in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>What is the spatial relationship of shape1 to shape2 in the image? Choices: A:{a} B:{b} C:{c} D:{d}. |
| Soild Geometry | cls | What is the shape of the {size} {color} object made of {material} in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Can you identify the type of {size} {color} object with {material} material in the picture? Choices: A:{a} B:{b} C:{c} D:{d}. |
| | cnt | There {be} {num} {shape} in the picture, right? Choices: A:Yes B:No.<br>{be} there {num} {shape} in the picture? Choices: A:Yes B:No.<br>There are {num} objects in the picture, right? Choices: A:Yes B:No.<br>Are there {num} shapes in the picture? Choices: A:Yes B:No.<br>How many {shape}s can you find in the picture? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Please count all the {shape}s in the image. Choices: A:{a} B:{b} C:{c} D:{d}.<br>Count the shapes in the image. Choices: A:{a} B:{b} C:{c} D:{d}.<br>How many shapes can you visually identify in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>How many {size} {color} {material} objects are there in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>How many {size} {color} {material} objects are present in the image? Choices: A:{a} B:{b} C:{c} D:{d}. |
| | grd | Please provide the bounding box coordinate of the region this sentence describes: {size} {color} {material} {shape}. |
| | rlat | Is it correct that {shape1} is described as being in the '{relation}' relative to {shape2} in the image? Choices: A:Yes B:No.<br>Is {shape1} described as being in the '{relation}' relative to {shape2} in the image? Choices: A:Yes B:No.<br>Is {shape1} said to be in the '{relation}' relative to {shape2} in the image? Choices: A:Yes B:No.<br>Can you confirm that {shape1} is in the '{relation}' relative to {shape2} in the image? Choices: A:Yes B:No.<br>In the image, where is {shape1} in relation to {shape2}? Choices: A:{a} B:{b} C:{c} D:{d}.<br>What is the relative position of {shape1} to {shape2} in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>What is the spatial relationship of {shape1} to {shape2} in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Describe how {shape1} is situated relative to {shape2} in the image. Choices: A:{a} B:{b} C:{c} D:{d}. |
| Graphs | cls | What type of chart is shown in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Identify the type of chart in the image. Choices: A:{a} B:{b} C:{c} D:{d}.<br>Which of the following best describes the chart graph in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Which kind of chart does the image represent? Choices: A:{a} B:{b} C:{c} D:{d}.<br>What kind of data visualization is shown in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Can you identify the type of chart in the image? Choices: A:{a} B:{b} C:{c} D:{d}. |
| | cnt | {be} there {num} {shape} in the picture? Choices: A:Yes B:No.<br>The picture contains {num} {shape}, each represented by a different color, doesn't it? Choices: A:Yes B:No.<br>You can see {num} {shape} with different colors in the picture, can't you? Choices: A:Yes B:No.<br>How many {shape}s are there in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>Can you count the total number of {shape}s in the image? Choices: A:{a} B:{b} C:{c} D:{d}.<br>What is the total number of {shape}s in the picture? Choices: A:{a} B:{b} C:{c} D:{d}. |
| | grd | Please provide the bounding box coordinate of the region this sentence describes: {color} {shape}. |
| | rlat | {question} Choices: A:Yes B:No. |

**Question:**
Find x.

**Choices:**
A: 2          B: $2\sqrt{11}$
C: 10        D: $4\sqrt{11}$

**SVE-Math-DeepSeek$^+$**

✔ Given that the length of the hypotenuse is 12, we can use the Pythagorean theorem to find the length of the perpendicular.

✔ Let's denote the length of the perpendicular as 'x'

✔ According to the Pythagorean theorem, we have $12^2 = 10^2 + x^2$, then $x = \sqrt{44} = 2\sqrt{11}$. Hence, the answer is option **B**
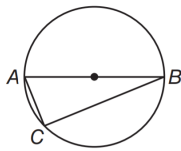
**SVE-Math-DeepSeek**

✘ Given that the area of triangle ABC is 10, we can use the formula for the area of a triangle to solve for x

✘ The area of triangle ABC is equal to the sum of the areas of triangle ACD and triangle BCD

✘ Since the height of triangle ACD is 4 and the height of triangle BCD is 2, we can write the equation as $0.5 \times 4 \times x + 0.5 \times 2 \times 5 = 10$. Then, we can find $x = 5$.

**InternVL2.5**

✘ To find x, we can use the Pythagorean theorem in the right triangle formed by x, y, and the hypotenuse z

✘ The hypotenuse z is the sum of the two segments, $10 + 12 = 22$. The right triangle has legs x and y.

✘ Using the Pythagorean theorem: $x^2 + y^2 = 22^2 = 484$. We need more information about y to solve.

(a)



**Question:**
AC = 8 inches, and BC = 15 inches. Find the radius of the circle.

**Choices:**
A: 7.5        B: 8
C: 8.5        D: 17

**SVE-Math-DeepSeek$^+$**

✔ Since ∠ACB = 90°, AB is the diameter.

✔ Since $AB = \sqrt{8^2 + 15^2}) = 17$, the radius of the circle is $\frac{1}{2} AB = 8.5$ inches. Therefore, the answer is **C**

**SVE-Math-DeepSeek**

✘ Connect OA and OB. Since AC and BC are tangents to circle O, we have OA perpendicular to AC and OB perpendicular to BC. Therefore, angle OAC = angle OBC = 90°.

✘ Also, ∠AOB = 2 ∠C = 36°. Hence, triangle AOB is an isosceles triangle.

✘ So, $OA = OB = AB = \sqrt{AC^2 + BC^2} = 17$. Therefore, the radius of circle O is 17. Therefore, the answer is **D**

**InternVL2.5**

✔ To find the radius of the circle, we can use the Pythagorean theorem. The given points form a right triangle with the diameter of the circle as the hypotenuse.

✔ Diameter $AB = \sqrt{AC^2 + BC^2} = \sqrt{8^2 + 15^2} = 17$. The radius r is half of the diameter. So, $r = \frac{1}{2} AB = 8.5$ inches. Correct option is **C**

(b)

Figure 15. Response comparisons between SVE-Math-DeepSeek$^+$, SVE-Math-DeepSeek, and InterVL2.5.

**SVE-Math-DeepSeek⁺**

✔ Since PA and PB are tangents to the circle O, we have ∠OAP = ∠OBP = 90°.

✔ Also, ∠AOB = 2 ∠AOC =130°.

✔ Therefore, ∠P = 360° - 90° - 90° - 130° = 50°. Thus, the answer is **C**

**SVE-Math-DeepSeek**

✔ Because PA and PB are tangents of circle O, A and B are the points of tangency. Therefore, ∠PBO = ∠PAO = 90°

✘ Since ∠ P = 65°, we have ∠BOP = 180° - 65° = 115°. Therefore, ∠ AOB = 360° - 90° - 90° - 115° = 65°. Thus, ∠P = 180° - ∠AOB = 115°. Therefore, the answer is **B**

**InternVL2.5**

✔ Tangents drawn from an external point to a circle are equal in length. Therefore, PA = PB. And ∠PAB= ∠PBA.

✘ ∠COB is given as 65°. Since O is the center of the circle, ∠COB is the central angle subtended by the arc AB.

✘ The inscribed ∠CAB is half of the central angle ∠COB. Therefore, ∠CAB = $\frac{65°}{2}$ = 32.5°.

✘ In triangle PAB, the sum of angles is 180°. So, ∠PAB + ∠PBA + ∠P = 180°. Then, ∠P = 180° − ∠PBA − ∠PAB = 180° − 65° = 115°. So, the answer is not in list.

**Question:**
As shown in the figure, PA and PB are tangent to circle O at A and B respectively. Then the degree of ∠P is ()?

**Choices:**
A: 65°      B: 135°
**C: 50°**      D: 100°

(a)



**SVE-Math-DeepSeek⁺**

✔ Since PA and PB are tangents to the circle O, we have ∠OAP = ∠OBP = 90°.

✔ Also, ∠P = 50°. Therefore, ∠AOB = 360° - 90° - 90° - 50° = 130°.

✔ Since ∠ACB = 1/2 ∠AOB (the central angle is half of the inscribed angle formed by the same arc), we have ∠ACB = 65°. Therefore, the correct answer is **A**

**SVE-Math-DeepSeek**

✔ Connect OA and OB. Since PA and PB is a tangent to circle O, we have ∠OAP = ∠OBP = 90°.

✔ Also, ∠P = 50°. Therefore, ∠AOB = 360° - 90° - 90° - 50° = 130°.

✔ Thus, ∠ACB = 1/2 ∠AOB = 65°. Therefore, the answer is **A**

**InternVL2.5**

✔ To find the size of angle ACB, we can use the property of the inscribed angle in a circle. The inscribed angle theorem states that an inscribed angle is half the measure of the central angle that subtends the same arc

✘ In the given figure, ∠ P is an exterior angle to triangle ABC, and it is given as 50°. According to the exterior angle theorem, the exterior angle is equal to the sum of the two opposite interior angles. Therefore, angle P is equal to the sum of ∠ACB and ∠ABC.

✘ Thus, ∠ACB = $\frac{1}{2}$ ∠P = $\frac{1}{2}$ × 50° = 25°. The answer is **C**

**Question:**
As shown in the figure, given that angle P = 50.0, then the size of angle ACB is?

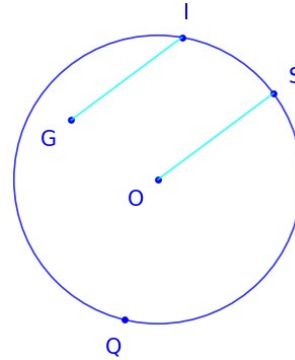**Choices:**
A: 65°      B: 60°
C: 55°      D: 50°

(b)

Figure 16. Response comparisons between SVE-Math-DeepSeek⁺, SVE-Math-DeepSeek, and InterVL2.5.

## Chain-of-Thought (CoT) Errors

**Question:** Can you explain how line GI and line OS in the image are related to each other?

*A.* parallel
*B.* perpendicular
*C.* on
*D.* incircle

*Ground Truth: A*

**GPT-4o Answer**
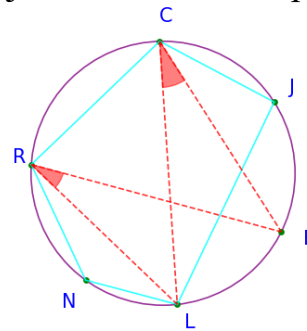
Step1: The lines GI and OS are both chords of the circle.
Step2: They are not parallel or perpendicular to each other, nor are they tangent or secant lines.
Step3: Therefore, the correct choice is C: on

## Chain-of-Thought (CoT) Errors

**Question:** Please ignore any shapes formed by red auxiliary lines. Can you identify the type of object CJLNR in the picture?

*A.* isosceles trapezoid
*B.* equilateral quadrilateral
*C.* pentagon
*D.* quadrangle

*Ground Truth: C*

**GPT-4o Answer**

Step1: The object labeled CJLNR in the picture is a quadrangle, which is a four-sided polygon.
Step2: The red auxiliary lines do not form any additional sides or angles that would change the classification of the shape.
Step3: Therefore, the correct answer to the question is D: quadrangle.

Figure 17. Examples of Chain-of-Thought (CoT) Errors in GPT-4o.

# Chain-of-Thought (CoT) Errors

**Question:** What is the shape of object IPRY in the image?

*A.* right trapezoid
*B.* isosceles trapezoid
*C.* equilateral quadrilateral
*D.* equilateral triangle

*Ground Truth: B*

**GPT-4o Answer**

Step1: To determine the shape of object IPRY, we need to analyze the angles and sides.

Step2: The angles at I and P appear to be right angles, and the sides IP and PR are parallel.

Step3: This suggests that IPRY is a right trapezoid. So, the answer is A: right trapezoid.

Figure 18. An example of Chain-of-Thought (CoT) Errors in GPT-4o.

# Chain-of-Thought (CoT) Errors

**Question:** Can you identify the type of object in the picture?

*A.* isosceles triangle
*B.* circle
*C.* scalene triangle
*D.* right triangle

*Ground Truth: D*

**Qwen2.5-VL-7B Answer**

Step1: The image shows a triangle with three vertices labeled as D, C, and E.

Step2: To determine the type of triangle, we need to look at the lengths of its sides.

Step3: If two sides are equal, it is an isosceles triangle. If all sides are of different lengths, it is a scalene triangle. If one angle is 90 degrees, it is a right triangle.

Step4: From the image, we can see that the lengths of the sides DE, EC, and DC are not equal, so the triangle is not isosceles or right.

Step5: Therefore, the correct answer is C: scalene triangle

Figure 19. An example of Chain-of-Thought (CoT) Errors in Qwen2.5-VL-7B.

# Chain-of-Thought (CoT) Errors

**Question:** What is the total number of shapes in the picture?

*A.* 3
*B.* 5
*C.* 4
*D.* 1



*Ground Truth: A*

**Qwen2.5-VL-7B Answer**

To determine the total number of shapes in the picture, let's analyze the image step by step:

Step1: The image shows a complex geometric figure with various lines and points.

Step2: The figure appears to be composed of multiple smaller shapes, but it is not immediately clear how many distinct shapes there are.

Step3: Upon closer inspection: There are several triangles formed by the lines connecting the points. There are also some quadrilaterals and possibly other polygons, but they are not clearly defined as separate shapes due to overlapping or shared sides.

Step4: Given the complexity and the overlapping nature of the lines, it is challenging to count each shape distinctly.

Step5: However, based on the visible structure, we can identify at least 5 distinct shapes (triangles and quadrilaterals).

Step6: Therefore, the total number of shapes in the picture is B: 5

Figure 20. An example of Chain-of-Thought (CoT) Errors in Qwen2.5-VL-7B.

# Chain-of-Thought (CoT) Errors

**Question:** Is it accurate to describe the relationship between circle O and triangle HYP in the image as \"incircle\"?

*A*. Yes
*B*. No

*Ground Truth: A*

**Qwen2.5-VL-7B Answer**

Step1: To determine if the relationship between circle O and triangle HYP is an \"incircle,\" we need to understand what an incircle is.

Step2: An incircle of a triangle is a circle that is tangent to all three sides of the triangle.

Step3: In the given image: Circle O appears to be inside triangle HYP. The center of the circle, O, seems to be equidistant from the sides of the triangle (HYP).

Step4: However, the points where the circle touches the sides of the triangle are not clearly marked or labeled in the image.

Step5: Given these observations, it is not possible to definitively say that circle O is the incircle of triangle HYP based solely on the visual information provided.

Step5: Therefore, the correct answer is B: No

Figure 21. An example of Chain-of-Thought (CoT) Errors in Qwen2.5-VL-7B.

## Recognition Errors

**Question:** How many objects are there in the image?

*A.* 1
*B.* 5
*C.* 3
*D.* 2

*Ground Truth: D*

**GPT-4o Answer**

C: 3

## Recognition Errors

**Question:** You can see 2 objects in the picture, can't you?

*A.* Yes
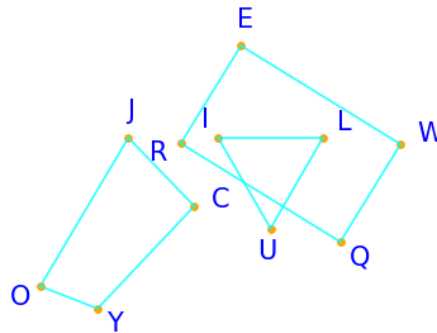*B.* No

*Ground Truth: B*

**GPT-4o Answer**

A: Yes

Figure 22. Examples of Recognition Errors in GPT-4o.
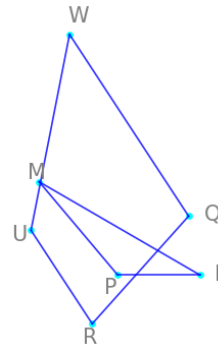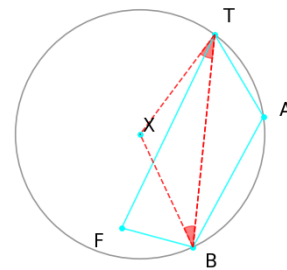
## Recognition Errors

**Question:** Can you identify the type of object QRUW in the picture?

*A.* isosceles trapezoid
*B.* trapezoid
*C.* pentagon
*D.* right trapezoid

*Ground Truth: B*

**GPT-4o Answer**

I can't identify the object QRUW in the picture.

## Recognition Errors

**Question:** Please provide the bounding box coordinate of the region this sentence describes: right triangle FGK.

*Ground Truth: [113, 182, 398, 329]*

**Qwen2.5-VL-7B Answer**

bbox_2d: [109, 135, 347, 286], label: right triangle FGK

Figure 23. Examples of Recognition Errors in GPT-4o (top) and Qwen2.5-VL-7B (bottom).

**Recognition Errors**

**Question:** What is the shape of object AFJP in the image?

*A.* isosceles triangle
*B.* equilateral quadrilateral
*C.* rectangle
*D.* right trapezoid

*Ground Truth: C*

**Qwen2.5-VL-7B Answer**

B



**Recognition Errors**

**Question:** Please ignore any shapes formed by red auxiliary lines. What is the shape of object ABTX in the image?

*A.* trapezoid
*B.* circle
*C.* square
*D.* ellipse

*Ground Truth: A*

**Qwen2.5-VL-7B Answer**

The correct answer is C: square

Figure 24. Examples of Recognition Errors in Qwen2.5-VL-7B.
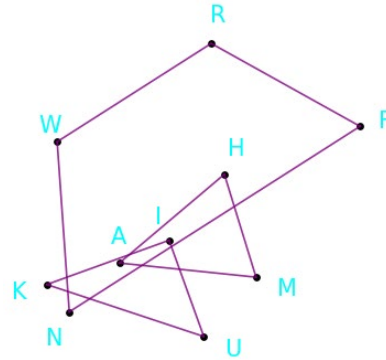
## Recognition Errors

**Question:** Can you identify the type of object AHM in the picture?

*A.* right triangle
*B.* isosceles triangle
*C.* rectangle
*D.* equilateral triangle

*Ground Truth: B*

**SVE-Math-DeepSeek⁺ Answer**

A

## Recognition Errors

**Question:** Please provide the bounding box coordinate of the region this sentence describes: quadrangle GJSY.

*Ground Truth: [243, 186, 343, 350]*

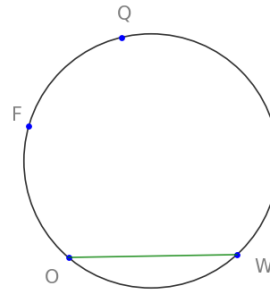**SVE-Math-DeepSeek⁺ Answer**

[287, 159, 399, 297]

Figure 25. Examples of Recognition Errors in SVE-Math-DeepSeek⁺-7B.

## Recognition Errors

**Question:** What is the shape of object IPRY in the image?

*A.* right trapezoid
*B.* isosceles trapezoid
*C.* equilateral quadrilateral
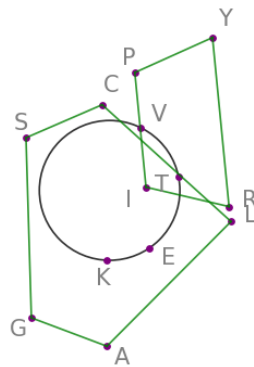*D.* equilateral triangle

*Ground Truth: B*

**SVE-Math-DeepSeek⁺ Answer**

C

Figure 26. An example of Recognition Errors in SVE-Math-DeepSeek $^+$-7B.