

# Plane Geometry Problem Solving with Multi-modal Reasoning: A Survey

Anonymous ACL submission

## Abstract

Plane geometry problem solving (PGPS) has recently gained significant attention as a benchmark to assess the multi-modal reasoning capabilities of large vision-language models. Despite the growing interest in PGPS, the research community still lacks a comprehensive overview that systematically synthesizes recent work in PGPS. To fill this gap, we present a survey of existing PGPS studies. We first categorize PGPS methods into an encoder-decoder framework and summarize the corresponding output formats used by their encoders and decoders. Subsequently, we classify and analyze these encoders and decoders according to their architectural designs. Finally, we outline major challenges and promising directions for future research. In particular, we discuss the hallucination issues arising during the encoding phase within encoder-decoder architectures, as well as the problem of data leakage in current PGPS benchmarks.

## 1 Introduction

Automated plane geometry problem solving (PGPS) has emerged as an important benchmark in artificial intelligence research due to its unique requirement for multi-modal reasoning with mathematical rigor (Seo et al., 2015; Chen et al., 2021). Typically, geometry problems combine textual descriptions with visual diagrams, each providing essential complementary information. The inherent necessity to integrate linguistic and visual modalities makes plane geometry a compelling testbed for advancing the multi-modal understanding capabilities of AI systems. Furthermore, practical motivations such as developing intelligent tutoring systems (Ritter et al., 2010; Alevn and Koedinger, 2002; Lee et al., 2025) and standardized benchmarks for evaluating AI reasoning (Chen et al., 2021; Cao and Xiao, 2022) highlight the importance of continued research in this area.

Nevertheless, substantial challenges persist in achieving full automation. Foremost among these is the complexity arising from the multi-modal nature of geometry problems, requiring precise alignment between textual statements and corresponding diagram elements (Seo et al., 2014). Resolving ambiguities in textual descriptions through visual references and accurately mapping entities between text and diagrams pose significant hurdles (Sachan et al., 2017; Zhang et al., 2022). Geometric diagrams also introduce unique challenges absent in natural images and other types of diagrams, including precise recognition of abstract symbols, e.g., angle markers and length indicators, accurate detection of geometric primitives, e.g., points, lines, and circles, and interpretation of implicit spatial relationships governed by geometric constraints. Additionally, effective PGPS demands embedding deep geometric domain knowledge, applying geometric axioms and theorems during the reasoning that are often implicitly assumed (Sachan et al., 2017; Sachan and Xing, 2017; Lu et al., 2021). Thus, integrating linguistic comprehension, visual analysis, and geometric reasoning continues to drive the complexity and significance of research in automated PGPS.

Recently, numerous new benchmarks, large-scale datasets, and model architectures have been proposed to tackle the challenges of PGPS. However, despite this rapid progress, most existing surveys on mathematical or multi-modal reasoning address geometry problems only as part of broader domains (Li et al., 2025; Yan et al., 2025; Yuan et al., 2025) and thus fail to examine the unique challenges of PGPS in depth. Consequently, the literature still lacks a dedicated, up-to-date survey centered on PGPS. The goal of this paper is to fill the gap by providing the PGPS research community with a structured overview of the latest benchmarks, datasets, and multi-modal reasoning approaches tailored specifically to PGPS.

The structure of this paper is summarized as follows: We first describe the definition of PGPS and relevant tasks (§2). We then introduce an overall framework for solving PGPS problems as an encoder-decoder architecture with intermediate representations (§3). Next, we review the details of encoder (§4) and decoder (§5) structures. Some additional thoughts are provided from the data collection perspective (Appendix A). Finally, we address the remaining challenges and promising future directions in automated PGPS (§6).

## 2 Tasks and benchmarks

In this section, we first define the PGPS and then introduce three tasks that are commonly tackled in the PGPS community, along with the benchmarks for each task.

### 2.1 Definition of PGPS

Euclidean plane geometry studies the properties and relationships among geometric primitives, e.g., points, lines, and circles, in a flat, two-dimensional space (Fitzpatrick and Heiberg, 2007). PGPS involves inferring unknown geometric properties or relationships from a given set of primitives and their known relations, such as determining the length of an unknown side in a triangle given the lengths of two sides and the measure of the included angle.

In real-world scenarios, plane geometry problems usually present as diagram and textual description pairs, as demonstrated in Fig. 1. The diagrams and accompanying textual descriptions typically complement each other in representing geometric primitives and relations. Diagrams usually provide visual information about spatial relationships, whereas textual descriptions explicitly mention properties or relational details. Due to this complementary nature, PGPS methods in real-world applications must not only infer unknown geometric facts but also accurately parse geometric information from these diagrams and text pairs.

### 2.2 PGPS tasks

We describe the three main tasks, along with the corresponding benchmarks, that are mainly tackled via PGPS research. Fig. 1 illustrates three examples for each task. For further details on the benchmarks from various perspectives, such as reasoning complexity, diagram-text interdependency, and data collection methods, refer to Appendix A.

#### 2.2.1 Direct-answer and multiple-choice tasks

**Task description** Most PGPS works quantify the capacity of a PGPS method to infer a single, well-defined property of a geometric entity from a unified diagrammatic-textual problem statement. The requested properties fall into two categories: i) numerical targets, e.g., angle magnitude, segment length, or area (Seo et al., 2015; Lu et al., 2021; Chen et al., 2021), and ii) categorical targets, e.g., the perpendicularity or parallelism of two lines (Xu et al., 2025).

PGPS methods are also evaluated through multiple-choice tasks (Lu et al., 2024; Zhang et al., 2025a). While these tasks use the same problems as direct-answer tasks, each multiple-choice problem provides a fixed set of candidate responses. A PGPS method must select the option that correctly identifies the target property, or equivalently, predict a value matching one of the provided choices. For example, in the scenario depicted in Fig. 1, the correct response is the label "c" or its corresponding value, "None."

**Evaluation metrics** In direct-answer tasks, performance is reported as top- $N$  accuracy: a PGPS method is considered correct when the ground truth answer appears within its  $N$  candidate answers. For multiple-choice tasks, the metric depends on the output representation of the method. If the method predicts an option label, evaluation reduces to standard top-1 accuracy. If it produces a value, e.g., scalar, a modified version of top- $N$  accuracy is utilized: the  $N$  generated values are scanned in order, and the attempt is scored correct once the first value that coincides with any listed option matches the ground truth.

**Benchmarks** Most PGPS benchmarks have been proposed to evaluate model performance on direct-answer and multiple-choice tasks. Some benchmarks exclusively consist of plane geometry problems (Alvin et al., 2017; Seo et al., 2015; Lu et al., 2021; Chen et al., 2021; Cao and Xiao, 2022; Zhang et al., 2023, 2024c; Fu et al., 2025; Kazemi et al., 2024; Xu et al., 2025), while others include plane geometry problems as part of broader benchmarks designed for general multi-modal reasoning evaluation (Lu et al., 2024; Zhang et al., 2025a; Yue et al., 2024; Kamoi et al., 2024; Wang et al., 2024a; Zou et al., 2025; Gupta et al., 2024; Wang et al., 2025).

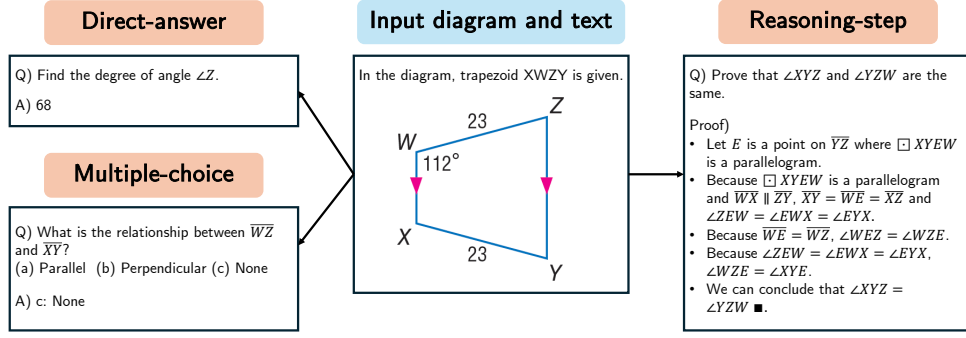


Figure 1: Illustration of three PGPS tasks. The three tasks are commonly used to evaluate PGPS methods in existing benchmarks: i) direct-answer, ii) multiple-choice, and iii) reasoning-step construction. In the direct-answer task, the model must predict a single numerical value as the answer to the problem. In the multiple-choice task, the model must select the correct label corresponding to the ground-truth option. In the reasoning-step construction task, the model is asked to generate the complete sequence of reasoning steps that lead to the correct final answer.

### 2.2.2 Reasoning tasks

**Task description** Some PGPS benchmarks assess methods not only on the correctness of the final answer but also on the soundness of the intermediate reasoning (Chen et al., 2022; Jaiswal et al., 2024). In a widely adopted proving problem setting, a PGPS method must generate a sequence of geometric axioms and theorems that derive the target statement, e.g., two angles are congruent, directly from the given conditions.

**Evaluation metrics** For reasoning-step construction tasks, top- $N$  accuracy is again adopted, granting success when any of the  $N$  predicted reasoning steps exactly reproduces the ground-truth steps.

**Benchmarks** UniGeo (Chen et al., 2022) is currently the only benchmark designed explicitly to systematically measure reasoning capabilities. Recently, approaches leveraging LLMs have emerged to evaluate individual reasoning steps (Zhang et al., 2025a; Jaiswal et al., 2024). However, these methods inherently rely on LLMs, posing significant limitations. Consequently, proposing diverse and systematic reasoning benchmarks remains an open research challenge.

## 3 Overall approach

PGPS models typically employ an encoder-decoder architecture, as demonstrated in Fig. 2. The *encoder* jointly processes the diagram and textual description to produce an *intermediate representation* that captures essential geometric information of the problem. The *decoder* then utilizes the extracted intermediate representation to generate a solution, presented as either a theorem sequence, a logic pro-

gram, or a natural-language description. Finally, the answer is obtained by post-processing the generated solution, e.g., by executing the logic program or extracting the final result from the natural-language description.

Before we discuss the detailed approach to constructing the encoder and decoder, we first review the output formats of the encoder and decoder commonly used across different PGPS tasks.

### 3.1 Encoder outputs

The output of an encoder forms an intermediate representation that can be further used as an input to a decoder. We categorize the output format of the encoder into i) formal-language description and ii) embedding vectors.

**Formal-language description** Several studies explicitly extract geometric primitives and relations from given diagram-text pairs, converting them into formal-language descriptions. A formal-language description consists of an *entity* set and a *predicate* set. The entity set contains geometric primitives, e.g., elementary primitives such as points, lines, and circles (Zhang et al., 2022, 2023), or higher-level shapes such as triangles and squares (Seo et al., 2015; Sachan et al., 2017; Sachan and Xing, 2017; Lu et al., 2021), along with non-geometric tokens such as numbers and variable names. The predicates define the relationships among the entities. For instance, an equality predicate binds two entities  $\angle ABC$  and  $30^\circ$  to represent the numerical value of the angle, i.e.,  $\angle ABC = 30^\circ$  or specify geometric relations, such as segments  $AB$  and  $BC$  being perpendicular, i.e.,  $AB \perp BC$ .

In earlier studies, rule-based approaches (Koo

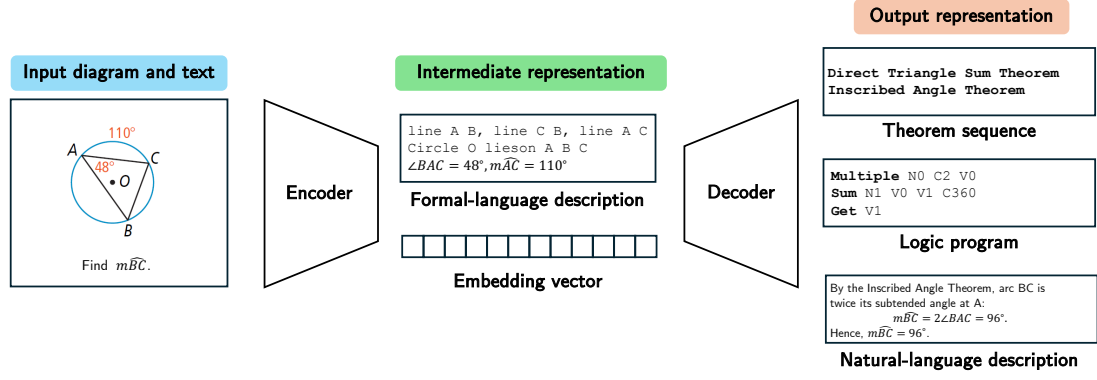


Figure 2: Visualization of the overall structure of PGPS methods. PGPS methods first encode the input diagram and text into an intermediate representation. The encoded representation is then passed to the decoder, which generates the final solution as a theorem sequence, a logic program, or a natural-language description.

et al., 2008; Bansal et al., 2014) and semantic parsers (Lewis et al., 2020) have been proposed to extract formal-language descriptions from textual descriptions without analyzing the diagram (Seo et al., 2015; Lu et al., 2021). Recent works extend these approaches to extract a formal language description from a diagram-text pair. Consequently, many PGPS studies release datasets consisting of diagrams and formal-language description pairs to train diagram parsers in a supervised way (Seo et al., 2015; Sachan et al., 2017; Sachan and Xing, 2017; Zhang et al., 2022; Lu et al., 2021; Zhang et al., 2023, 2024c)

**Embedding vectors** Certain PGPS encoders represent inputs as embedding vectors, typically utilizing one of three strategies: i) embedding diagrams and textual descriptions separately and subsequently merging them (Chen et al., 2021; Cao and Xiao, 2022; Chen et al., 2022; Ning et al., 2023; Liang et al., 2023; Jian et al., 2023b), ii) embedding diagrams exclusively and then combining them with raw textual inputs (Xia et al., 2025; Cho et al., 2025; Shi et al., 2024; Zhang et al., 2025b; Gao et al., 2025; Zhang et al., 2025e; Peng et al., 2025; Xu et al., 2024), or iii) jointly processing diagrams and texts through a unified encoder (Zhang et al., 2023; Li et al., 2024). Although these embeddings are generally less interpretable compared to formal-language descriptions, they enable end-to-end training with the decoder.

### 3.2 Decoder outputs

Given the output of the encoder, the decoder generates the solution from which the final answer can be derived. We classify decoder output formats into three types: i) theorem sequences, ii) logic programs, and iii) natural-language descriptions.

**A sequence of theorems** Many PGPS works represent the output of a PGPS problem as a sequence of theorem applications. This approach naturally aligns with a reasoning process, in which theorems are iteratively applied to given entities and predicates to logically derive new geometric facts, including the target predicate specified as a goal (Trinh et al., 2024). Specifically, given geometric entities and predicates extracted from the original description, theorems from a predefined library can be applied to the entities and predicates to derive additional predicates not explicitly stated in the original problem. Recent PGPS datasets provide annotated triples consisting of the formal-language description, the target predicate, and a corresponding reference theorem sequence (Sachan et al., 2017; Sachan and Xing, 2017; Lu et al., 2021; Zhang et al., 2024c).

**A logic program** A logic program is commonly adopted as an output representation for PGPS. Specifically, inspired by the observation that the reasoning process in PGPS typically involves applying a series of operations to numerical constants and variables provided in the problem (Chen et al., 2021; Amini et al., 2019; Chen et al., 2023), a logic program is defined as a sequence of triples, each consisting of an operation and its operands, such as numerical values and variable names. The operations in these programs fall into two main categories: i) arithmetic functions, ranging from basic operations like addition and multiplication to geometry-specific computations such as the Pythagorean operation (Chen et al., 2021; Cao and Xiao, 2022; Chen et al., 2022), and ii) equality assertions that establish identity between two expressions (Zhang et al., 2023). Several PGPS datasets provide paired examples, each consisting



of a diagram-text problem and its corresponding logic program (Chen et al., 2021; Cao and Xiao, 2022; Chen et al., 2022; Zhang et al., 2023).

**A natural-language description** Recent PGPS methods generate solutions and answers in natural language without relying on a specific template. The inherent flexibility of natural language allows these models to easily provide outputs for a wide range of tasks, e.g., geometric diagram captioning, without being limited to fixed problem-solving formats. To train such methods, various types of PGPS datasets have been proposed. For tasks which focus on problem solving, the output, given a diagram and text, can either be the answer expressed in natural language (Shi et al., 2024) or a reasoning path in the form of a chain-of-thought (Wei et al., 2022) to infer the answer (Zhang et al., 2025b; Gao et al., 2025). In addition to problem solving, datasets have also been proposed for tasks such as geometric diagram captioning (Zhang et al., 2025b; Gao et al., 2025; Cho et al., 2025; Xia et al., 2025) and question answering (Gao et al., 2025).

### 3.3 Encoder-decoder with desired outputs

Once the intermediate representations and output representations are determined based on target problems or tasks, one can choose an appropriate encoder and decoder that can produce the desired outputs. Fig. 3 summarizes possible combinations of encoder-decoder architectures along with the desired outputs. A combination of encoder, intermediate representation, decoder, and output representation can lead to a specific architecture for PGPS. In the following two sections, we review the possible choices of encoder and decoder structures.

## 4 Encoders

The encoder extracts the relevant components from the given diagram and text that are necessary for PGPS. We review the encoders in the following aspects: i) rule-based and ii) neural network-based.

### 4.1 Rule-based encoders

Early PGPS methods relied on classical computer vision and text parsing techniques to independently extract geometric primitives and relations from diagrams and text, merging them into formal-language descriptions. Most studies (Seo et al., 2015; Sachan et al., 2017; Sachan and Xing, 2017; Alvin et al., 2017; Gan et al., 2019) employed rule-based diagram parsers, notably HoughGeo (Chen et al.,

2015) or G-Aligner (Seo et al., 2014), which preprocess diagrams to detect geometric primitives using classical computer vision techniques, e.g., Gaussian blur and Hough transforms, and then match detected primitives to literal sets using either handcrafted rules or optimization. For textual extraction, many approaches (Wu et al., 2024b; Zhao et al., 2025; Peng et al., 2023; Zhang et al., 2024b; Jian et al., 2023a; Zou et al., 2024) adopted the InterGPS (Lu et al., 2021) parser, a rule-based method utilizing regular expressions, which is reliable and effective even with limited data.

### 4.2 Neural network encoders

We review the neural network-based encoders based on the desired output format.

#### 4.2.1 Formal-language description generation

Recent PGPS approaches adopt neural encoders to generate formal-language descriptions from diverse diagrams and texts, typically training separate encoders for each modality. Neural diagram encoders commonly operate in two stages: primitive detection using object detectors such as RetinaNet (Lin et al., 2017b; Lu et al., 2021) and feature pyramid networks (Lin et al., 2017a; Zhang et al., 2022), followed by relation inference modeled either as a constrained optimization problem (Lu et al., 2021) or as a graph-learning task leveraging graph neural networks (GNNs) (Zhang et al., 2022). For text encoding, subsequent PGPS studies (Sachan et al., 2017; Sachan and Xing, 2017) commonly employ logistic regression models, as originally introduced by GEOS (Seo et al., 2015), to extract primitives and relations from problem statements.

#### 4.2.2 Embedding vector generation

To enable end-to-end learning, recent PGPS methods employ neural encoders that map both the diagram and text into a unified embedding space, providing a joint vector representation for PGPS. Here, we review the neural encoders based on their training strategy.

**Learning from scratch** Early PGPS works train joint diagram-text encoders and decoders end-to-end from scratch on target PGPS datasets. Diagram embeddings commonly utilize convolutional neural networks (CNNs), including vanilla CNN (Zhang et al., 2023), ResNet (He et al., 2016; Chen et al., 2021; Cao and Xiao, 2022), DenseNet (Huang et al., 2017; Jian et al., 2023a), and VQ-VAE encoders (van den Oord et al., 2017; Liang et al.,

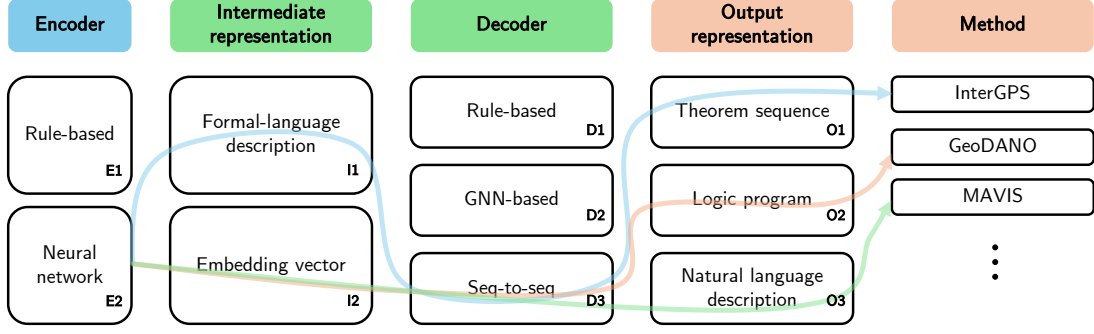


Figure 3: Overview of the PGPS pipeline. PGPS methods can be categorized based on the combination of the encoder, intermediate representation, decoder, and output representation. For example, the InterGPS can be represented as a combination of E2, I1, D3, and O1. We summarize PGPS methods as a combination of these components in Table A1.

2023), as well as Vision Transformers (ViT) (Dosovitskiy et al., 2021; Ning et al., 2023). Text embeddings are typically produced by sequential models like LSTMs (Hochreiter and Schmidhuber, 1997; Chen et al., 2021; Cao and Xiao, 2022) or Transformer-based encoders (Vaswani et al., 2017), such as vanilla Transformer (Zhang et al., 2023; Ning et al., 2023; Li et al., 2024) and RoBERTa (Liu et al., 2019; Cao and Xiao, 2022). Diagram and text embeddings are fused via co-attention networks (Yu et al., 2019; Chen et al., 2021; Ning et al., 2023), bi-directional GRUs (Chung et al., 2014; Zhang et al., 2023; Li et al., 2024), or Transformers (Chen et al., 2022).

Besides direct optimization on PGPS tasks, joint encoders frequently employ auxiliary objectives for improved performance. Many approaches incorporate self-supervised tasks, including jigsaw-location prediction (Chen et al., 2021; Cao and Xiao, 2022; Jian et al., 2023a), masked-token prediction in text (Devlin et al., 2019; Chen et al., 2022; Zhang et al., 2023; Li et al., 2024) or diagrams (He et al., 2022; Ning et al., 2023), text-conditioned diagram-symbol classification (Ning et al., 2023), and VQ-VAE objective (Liang et al., 2023). Other studies leverage explicit labels, training encoders for geometry-element or knowledge-point classification (Chen et al., 2021; Cao and Xiao, 2022), or contrastive learning between diagram patches and textual tokens (Li et al., 2024).

**Pre-trained encoders** To leverage pretrained knowledge and enhance training efficiency, many recent PGPS methods employ neural encoders inspired by the LLaVA architecture (Liu et al., 2023), which integrates a pretrained vision encoder to encode diagrams. Specifically, diagrams are first

transformed into visual embeddings using a pre-trained vision encoder, followed by a lightweight adapter consisting of a multi-layer perceptron. During training, only the adapter parameters are updated, keeping the vision encoder frozen to preserve general visual knowledge and reduce training cost. While OpenCLIP (Radford et al., 2021) is the most commonly used backbone (Shi et al., 2024; Gao et al., 2025; Xu et al., 2024), other general-purpose models such as SigLIP (Zhai et al., 2023; Zhang et al., 2025e) and InternViT (Chen et al., 2024c; Peng et al., 2025), as well as the math-specific Math-CLIP encoder (Zhang et al., 2025b; Peng et al., 2025), have also been employed.

**Fine-tuned encoders** Most pretrained vision encoders perform poorly when applied to geometric diagrams (Zhang et al., 2025b; Xia et al., 2025; Cho et al., 2025). To address this limitation, PGPS methods employing the LLaVA-style architecture typically fine-tune the vision encoders before integrating them into downstream pipelines. Two main fine-tuning strategies are common: i) self-supervised methods such as masked auto-encoding (He et al., 2022; Xia et al., 2025), and ii) weakly supervised methods such as CLIP (Zhang et al., 2025b; Cho et al., 2025), direct preference optimization (Rafailov et al., 2023; Huang et al., 2025), or grounding tasks (Li\* et al., 2022; Zhang et al., 2025c), which leverage synthetic geometric diagrams and labels pairs. Nevertheless, since synthetic diagrams do not fully capture the characteristics of real-world diagrams, GeoDANO (Cho et al., 2025) further employs few-shot domain adaptation under the same CLIP training objective to minimize the residual domain gap.

## 5 Decoders

Based on the representations produced by the encoder, the decoder generates the solution to the problem. We survey the PGPS decoders using the following dimensions: i) input representation and ii) architectural design.

### 5.1 Formal-language description decoder

We first introduce the architectures of the decoders that receive a formal language description as input.

**Rule-based axiomatic decoders** Several methods that operate on formal-language descriptions determine the required theorem sequence with a rule-based decoder. GEOS++ (Sachan et al., 2017) employs an exhaustive brute-force search to locate a sequence of theorems whose application yields the target predicate. GeoShader (Alvin et al., 2017) specifies a deterministic set of composition rules that directly selects the relevant theorems without search. GEOS-OS (Sachan and Xing, 2017) trains a log-linear model to assign probabilities to candidate theorems and then performs beam search, returning the highest-scoring theorem sequence.

**GNN-based decoders** A formal-language description, composed of geometric primitives and their relations, naturally corresponds to a graph structure. Exploiting this, several PGPS decoders first encode the formal description as a graph or hypergraph and then generate theorem-application sequences from the resulting graph representation. Such encodings typically follow one of three schemes: i) primitives as nodes and predicates as edges (Peng et al., 2023), ii) primitives and predicates both as nodes connected via edges (Jian et al., 2023a), or iii) predicates as hypernodes and theorems as directed hyperedges forming a hypertree (Zhang et al., 2024b). These encoded structures are subsequently fed into graph-to-sequence decoders, such as Graphormer (Zhang et al., 2024b), graph Transformer (Peng et al., 2023), or graph convolutional network (Kipf and Welling, 2017) followed by LSTM (Jian et al., 2023a), to produce the target theorem sequence.

**Sequence-to-sequence decoders** Some approaches treat formal-language descriptions as a flat token sequence and pass it directly to a sequence-to-sequence (seq-to-seq) model to generate the corresponding theorem sequence. Transformers are predominantly employed for

these tasks by encoding the formal description directly (Lu et al., 2021; Wu et al., 2024b; Zou et al., 2024). A few studies instead utilize off-the-shelf LLMs, e.g., o3-mini (OpenAI, 2025b), without additional training (Zhao et al., 2025).

### 5.2 Seq-to-seq embedding decoders

Several PGPS studies feed either a joint diagram-text embedding or a concatenation of diagram embedding and raw text into a sequence-to-sequence decoder. Early work primarily employs RNN-based decoders such as LSTMs or GRUs (Chen et al., 2021; Cao and Xiao, 2022; Zhang et al., 2023; Li et al., 2024; Ning et al., 2023; Jian et al., 2023b), while later studies commonly adopt encoder-decoder Transformers such as T5 (Raffel et al., 2020; Liang et al., 2023; Chen et al., 2022). The recent proliferation of LLMs has motivated a shift toward fine-tuning encoder-only Transformers, such as LLaMA (Touvron et al., 2023; Cho et al., 2025; Gao et al., 2025; Xu et al., 2024) and Vicuna (Vicuna, 2023; Shi et al., 2024), specifically adapted for PGPS tasks.

## 6 Challenges and future directions

We examine the remaining challenges in PGPS and propose potential directions for future work.

### 6.1 Hallucination in diagram perception

PGPS methods initially extract geometric primitives and relations from diagrams and text, making accurate perception crucial before reasoning. However, studies indicate that PGPS methods frequently misperceive these primitives and relations, especially when generating natural-language descriptions (Huang et al., 2025; Zhang et al., 2025a) as depicted in Fig. A1. For example, Table A2 reveals that GPT-4.1 (OpenAI, 2025a) fails to capture a fundamental geometric relation among the points and lines and produces hallucinations. These hallucinations not only degrade PGPS performance but also diminish dataset quality. Computer vision studies report similar hallucination issues in datasets produced by large VLMs (Zhang et al., 2025d; Sahoo et al., 2024; Li et al., 2023; Chen et al., 2024b), further evidenced in PGPS datasets as shown in Table 1. Consequently, models trained on hallucinated data suffer measurable performance declines (Zhang et al., 2025d; Lai et al., 2025; Yu et al., 2024; Hirota et al., 2024).

Visual prompting techniques, such as augmenting diagrams with bounding boxes, markers, or

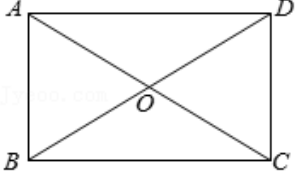
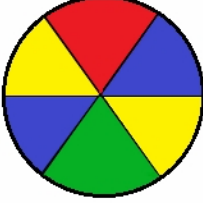
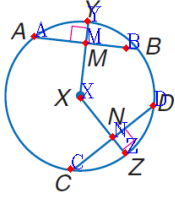
	Example 1	Example 2	Example 3
<b>Diagram</b>			
<b>Question</b>	In the given figure, let's denote the area of triangle AOB as variable $x$ . Find the area of rectangle ABCD in terms of $x$ . <b>Choices: A: 8 B: 10 C: 12 D: 16</b>	Based on the image, what is the measure of the interior angle at <b>vertex A</b> ? Choices: A. 90 degrees B. More than 90 degrees C. Less than 90 degrees D. Cannot be determined	Does the diagram include any line segments that are not perpendicular to each other?
<b>Solution</b>	To determine the area of rectangle ABCD, we can use the fact that triangle AOB is half the area of the rectangle. <b>Therefore, the area of rectangle ABCD is 2 times the area of triangle AOB, which is <math>2x</math>.</b> Hence, the answer is <b>option B. Answer:D</b>	Use the properties of the geometric shapes and theorems related to angles to deduce the measure of <b>the interior angle at vertex A</b> based on the given image and information. So the <b>answer is B</b>	Yes, in the diagram, <b>line segment YM is not perpendicular to line segment MA.</b>

Table 1: Examples of hallucinations in the natural-language description datasets annotated with L(V)LM. We visualize the examples from the PGPS datasets, e.g., G-LLaVA and MAVIS, which contain hallucinations in the question or response due to the L(V)LM annotation. We highlight the hallucinations with bold characters.

segmentation masks, have emerged as promising solutions for mitigating hallucinations (Wu et al., 2024a; Yang et al., 2023; Ma et al., 2025). These methods are especially beneficial for PGPS tasks, as they dynamically highlight relevant primitives and relations during reasoning and facilitate the critical step of drawing auxiliary lines. Augmenting diagrams at test time (Muennighoff et al., 2025) by applying segmentation masks (Ravi et al., 2024) or adding auxiliary constructions aligned with the current reasoning step (Murphy et al., 2024; Hu et al., 2024b) offers a practical approach to enhance multi-modal reasoning performance in PGPS.

## 6.2 Evaluation challenges in benchmarks

Comprehensive PGPS benchmarks should evaluate perception across diverse, realistic diagrams, ensuring that visual processing is essential for solving each problem. However, as shown in Table A3, existing benchmarks do not satisfy these criteria simultaneously. Synthetic diagrams, while scalable, often fail to represent the complexity of real-world scenarios (Zhong et al., 2025; Bates et al., 2025; Wang et al., 2024b), lacking elements such as parallel markers or placeholder objects, as illustrated in Fig. A2. Conversely, manually collected benchmarks better reflect real-world complexity

but frequently reuse diagrams from popular PGPS datasets, introducing data leakage and compromising domain generalization evaluations (Hu et al., 2024a; Cao et al., 2024; Chen et al., 2024a).

Even manually curated benchmarks without common PGPS dataset reuse often neglect crucial diagram-text dependencies discussed in Appendix A.2. MathVerse addresses these dependencies explicitly and avoids synthetic diagrams, but still suffers from data leakage, limiting its capability to assess genuine multi-modal reasoning. To overcome these issues, future research should develop synthetic diagram generators that closely replicate real-world complexity or create new datasets that strictly require visual reasoning while rigorously preventing data leakage.

## 7 Conclusion

In this paper, we examine the tasks, benchmarks, and methods used in existing PGPS research. We summarize the main PGPS approaches as an encoder-decoder architecture, along with the intermediate and output representations utilized across different methods. Through the analysis, we outline future research directions addressing current challenges, particularly regarding diagram perception and benchmark comprehensiveness.



## Limitations

In this paper, we primarily survey studies related to PGPS. While our work offers a comprehensive review of the existing PGPS literature, it is limited to two-dimensional geometry. Consequently, we do not address research involving three-dimensional geometry, such as projective and solid geometry, which requires understanding spatial relationships in three-dimensional space.

## References

- Vincent A.W.M.M. Aleven and Kenneth R. Koedinger. 2002. [An effective metacognitive strategy: learning by doing and explaining with a computer-based cognitive tutor](#). *Cognitive Science*, 26(2):147–179.
- Chris Alvin, Sumit Gulwani, Rupak Majumdar, and Supratik Mukhopadhyay. 2017. Synthesis of problems for shaded area geometry reasoning. In *Artificial Intelligence in Education*, pages 455–458, Cham. Springer International Publishing.
- Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. [MathQA: Towards interpretable math word problem solving with operation-based formalisms](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2357–2367, Minneapolis, Minnesota. Association for Computational Linguistics.
- Mohit Bansal, Kevin Gimpel, and Karen Livescu. 2014. [Tailoring continuous word representations for dependency parsing](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 809–815, Baltimore, Maryland. Association for Computational Linguistics.
- Averi Bates, Ryan Vavricka, Shane Carleton, Ruosi Shao, and Chongle Pan. 2025. [Unified modeling language code generation from diagram images using multimodal large language models](#). *Machine Learning with Applications*, 20:100660.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. [Language models are few-shot learners](#). *Preprint*, arXiv:2005.14165.
- Jie Cao and Jing Xiao. 2022. [An augmented benchmark dataset for geometric question answering through dual parallel text encoding](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 1511–1520, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Lele Cao, Valentin Buchner, Zineb Senane, and Fangkai Yang. 2024. [Introducing GenCeption for multimodal LLM benchmarking: You may bypass annotations](#). In *Proceedings of the 4th Workshop on Trustworthy Natural Language Processing (TrustNLP 2024)*, pages 196–201, Mexico City, Mexico. Association for Computational Linguistics.
- Jiaqi Chen, Tong Li, Jinghui Qin, Pan Lu, Liang Lin, Chongyu Chen, and Xiaodan Liang. 2022. [UniGeo: Unifying geometry logical reasoning via reformulating mathematical expression](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3313–3323, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Jiaqi Chen, Jianheng Tang, Jinghui Qin, Xiaodan Liang, Lingbo Liu, Eric Xing, and Liang Lin. 2021. [GeoQA: A geometric question answering benchmark towards multimodal numerical reasoning](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 513–523, Online. Association for Computational Linguistics.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, and 1 others. 2024a. Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*.
- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. 2023. [Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks](#). *Transactions on Machine Learning Research*.
- Xiaoyu Chen, Dan Song, and Dongming Wang. 2015. [Automated generation of geometric theorems from images of diagrams](#). *Annals of Mathematics and Artificial Intelligence*, 74(3):333–358.
- Xuwei Chen, Ziqiao Ma, Xuejun Zhang, Sihan Xu, Jianing Yang, David F. Fouhey, Joyce Chai, and Shengyi Qian. 2024b. [Multi-object hallucination in vision language models](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 44393–44418. Curran Associates, Inc.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. 2024c. 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 (CVPR)*, pages 24185–24198.
- Seunghyuk Cho, Zhenyue Qin, Yang Liu, Youngbin Choi, Seungbeom Lee, and Dongwoo Kim. 2025. [Geodano: Geometric vlm with domain agnostic vision encoder](#). *Preprint*, arXiv:2502.11360.

- Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. [Empirical evaluation of gated recurrent neural networks on sequence modeling](#). *Preprint*, arXiv:1412.3555.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. [An image is worth 16x16 words: Transformers for image recognition at scale](#). In *International Conference on Learning Representations*.
- R. Fitzpatrick and J.L. Heiberg. 2007. [Euclid’s Elements](#). University of Texas at Austin, Institute for Fusion Studies Department of Physics.
- Daocheng Fu, Zijun Chen, Renqiu Xia, Qi Liu, Yuan Feng, Hongbin Zhou, Renrui Zhang, Shiyang Feng, Peng Gao, Junchi Yan, Botian Shi, Bo Zhang, and Yu Qiao. 2025. [Trustgeogen: Scalable and formal-verified data engine for trustworthy multi-modal geometric problem solving](#). *Preprint*, arXiv:2504.15780.
- Wenbin Gan, Xinguo Yu, Ting Zhang, and Mingshu Wang. 2019. [Automatically proving plane geometry theorems stated by text and diagram](#). *International Journal of Pattern Recognition and Artificial Intelligence*, 33(07):1940003.
- Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, Lanqing HONG, Jianhua Han, Hang Xu, Zhenguo Li, and Lingpeng Kong. 2025. [G-LLaVA: Solving geometric problem with multi-modal large language model](#). In *The Thirteenth International Conference on Learning Representations*.
- Himanshu Gupta, Shreyas Verma, Ujjwala Ananthaswaran, Kevin Scaria, Mihir Parmar, Swaroop Mishra, and Chitta Baral. 2024. [Polymath: A challenging multi-modal mathematical reasoning benchmark](#). *Preprint*, arXiv:2410.14702.
- Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2022. [Masked autoencoders are scalable vision learners](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16000–16009.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. [Deep residual learning for image recognition](#). In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yusuke Hirota, Ryo Hachiuma, Chao-Han Huck Yang, and Yuta Nakashima. 2024. [From descriptive richness to bias: Unveiling the dark side of generative image caption enrichment](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 17807–17816, Miami, Florida, USA. Association for Computational Linguistics.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long short-term memory](#). *Neural Computation*, 9(8):1735–1780.
- Xuhao Hu, Dongrui Liu, Hao Li, Xuanjing Huang, and Jing Shao. 2024a. [Vlsbench: Unveiling visual leakage in multimodal safety](#). *arXiv preprint arXiv:2411.19939*.
- Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Ranjay Krishna. 2024b. [Visual sketchpad: Sketching as a visual chain of thought for multimodal language models](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. 2017. [Densely connected convolutional networks](#). In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Kung-Hsiang Huang, Can Qin, Haoyi Qiu, Philippe Laban, Shafiq Joty, Caiming Xiong, and Chien-Sheng Wu. 2025. [Why vision language models struggle with visual arithmetic? towards enhanced chart and geometry understanding](#). *Preprint*, arXiv:2502.11492.
- Raj Jaiswal, Avinash Anand, and Rajiv Ratn Shah. 2024. [Advancing multimodal llms: A focus on geometry problem solving reasoning and sequential scoring](#). In *Proceedings of the 6th ACM International Conference on Multimedia in Asia, MMAsia ’24*, New York, NY, USA. Association for Computing Machinery.
- Pengpeng Jian, Fucheng Guo, Cong Pan, Yanli Wang, Yangrui Yang, and Yang Li. 2023a. [Interpretable geometry problem solving using improved retinanet and graph convolutional network](#). *Electronics*, 12(22).
- Pengpeng Jian, Fucheng Guo, Yanli Wang, and Yang Li. 2023b. [Solving geometry problems via feature learning and contrastive learning of multimodal data](#). *Computer Modeling in Engineering & Sciences*, 136(2):1707–1728.
- Ryo Kamoi, Yusen Zhang, Sarkar Snigdha Sarathi Das, Ranran Haoran Zhang, and Rui Zhang. 2024. [Visonlyqa: Large vision language models still struggle with visual perception of geometric information](#). *arXiv preprint arXiv:2412.00947*.
- Mehran Kazemi, Hamidreza Alvari, Ankit Anand, Jialin Wu, Xi Chen, and Radu Soricut. 2024. [Geomverse: A systematic evaluation of large models for geometric reasoning](#). In *AI for Math Workshop @ ICML 2024*.





974	Maizhen Ning, Qiu-Feng Wang, Kaizhu Huang, and Xiaowei Huang. 2023. <a href="#">A symbolic characters aware model for solving geometry problems</a> . In <i>Proceedings of the 31st ACM International Conference on Multimedia</i> , MM '23, page 7767–7775, New York, NY, USA. Association for Computing Machinery.	1030
975		1031
976		1032
977		
978		1033
979		1034
980	OpenAI. 2023. Gpt-4v(ision) system card. <a href="https://openai.com/index/gpt-4v-system-card">https://openai.com/index/gpt-4v-system-card</a> .	1035
981		1036
982	OpenAI. 2025a. Introducing gpt-4.1 in the api. <a href="https://openai.com/index/gpt-4-1/">https://openai.com/index/gpt-4-1/</a> .	1037
983		1038
984	OpenAI. 2025b. Openai o3-mini. <a href="https://openai.com/index/openai-o3-mini">https://openai.com/index/openai-o3-mini</a> .	1039
985		
986	Shuai Peng, Di Fu, Yijun Liang, Liangcai Gao, and Zhi Tang. 2023. <a href="#">GeoDRL: A self-learning framework for geometry problem solving using reinforcement learning in deductive reasoning</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 13468–13480, Toronto, Canada. Association for Computational Linguistics.	1040
987		1041
988		1042
989		1043
990		1044
991		1045
992		
993	Tianshuo Peng, Mingsheng Li, Hongbin Zhou, Renqiu Xia, Renrui Zhang, Lei Bai, Song Mao, Bin Wang, Conghui He, Aojun Zhou, Botian Shi, Tao Chen, Bo Zhang, and Xiangyu Yue. 2025. <a href="#">Chimera: Improving generalist model with domain-specific experts</a> . <i>Preprint</i> , arXiv:2412.05983.	1046
994		1047
995		1048
996		1049
997		1050
998		1051
999		1052
1000	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. <a href="#">Learning transferable visual models from natural language supervision</a> . In <i>Proceedings of the 38th International Conference on Machine Learning</i> , volume 139 of <i>Proceedings of Machine Learning Research</i> , pages 8748–8763. PMLR.	1053
1001		
1002		
1003		
1004		
1005		
1006		
1007		
1008	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. <a href="#">Direct preference optimization: Your language model is secretly a reward model</a> . In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1054
1009		1055
1010		1056
1011		1057
1012		1058
1013		
1014	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>J. Mach. Learn. Res.</i> , 21(1).	1059
1015		1060
1016		1061
1017		1062
1018		1063
1019	Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. 2024. <a href="#">Sam 2: Segment anything in images and videos</a> . <i>arXiv preprint arXiv:2408.00714</i> .	1064
1020		1065
1021		1066
1022		1067
1023		1068
1024		1069
1025		1070
1026		1071
1027	Steven Ritter, Brendon Towle, R. Charles Murray, Robert G M. Hausmann, and John Connelly. 2010. <a href="#">A cognitive tutor for geometric proof</a> . In <i>Proceedings of the 10th International Conference on Intelligent Tutoring Systems - Volume Part II, ITS'10</i> , page 453, Berlin, Heidelberg. Springer-Verlag.	1072
1028		1073
1029		
	Mrinmaya Sachan, Kumar Dubey, and Eric Xing. 2017. <a href="#">From textbooks to knowledge: A case study in harvesting axiomatic knowledge from textbooks to solve geometry problems</a> . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 773–784, Copenhagen, Denmark. Association for Computational Linguistics.	1074
		1075
		1076
		1077
		1078
		1079
		1080
	Mrinmaya Sachan and Eric Xing. 2017. <a href="#">Learning to solve geometry problems from natural language demonstrations in textbooks</a> . In <i>Proceedings of the 6th Joint Conference on Lexical and Computational Semantics (*SEM 2017)</i> , pages 251–261, Vancouver, Canada. Association for Computational Linguistics.	1081
		1082
		1083
		1084
		1085
		1086
		1087
	Pranab Sahoo, Prabhaskar Meheria, Akash Ghosh, Sriparna Saha, Vinija Jain, and Aman Chadha. 2024. <a href="#">A comprehensive survey of hallucination in large language, image, video and audio foundation models</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 11709–11724, Miami, Florida, USA. Association for Computational Linguistics.	1088
		1089
		1090
		1091
		1092
	Min Joon Seo, Hannaneh Hajishirzi, Ali Farhadi, and Oren Etzioni. 2014. Diagram understanding in geometry questions. In <i>Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence, AAAI'14</i> , page 2831–2838. AAAI Press.	1093
		1094
		1095
		1096
		1097
		1098
		1099
		1100
		1101
		1102
		1103
		1104
		1105
		1106
		1107
		1108
		1109
		1110
		1111
		1112
		1113
		1114
		1115
		1116
		1117
		1118
		1119
		1120
		1121
		1122
		1123
		1124
		1125
		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145
		1146
		1147
		1148
		1149
		1150
		1151
		1152
		1153
		1154
		1155
		1156
		1157
		1158
		1159
		1160
		1161
		1162
		1163
		1164
		1165
		1166
		1167
		1168
		1169
		1170
		1171
		1172
		1173
		1174
		1175
		1176
		1177
		1178
		1179
		1180
		1181
		1182
		1183
		1184
		1185
		1186
		1187
		1188
		1189
		1190
		1191
		1192
		1193
		1194
		1195
		1196
		1197
		1198
		1199
		1200
		1201
		1202
		1203
		1204
		1205
		1206
		1207
		1208
		1209
		1210
		1211
		1212
		1213
		1214
		1215
		1216
		1217
		1218
		1219
		1220
		1221
		1222
		1223
		1224
		1225
		1226
		1227
		1228
		1229
		1230
		1231
		1232
		1233
		1234
		1235
		1236
		1237
		1238
		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246
		1247
		1248
		1249
		1250
		1251
		1252
		1253
		1254
		1255
		1256
		1257
		1258
		1259
		1260
		1261
		1262
		1263
		1264
		1265
		1266
		1267
		1268
		1269
		1270
		1271
		1272
		1273
		1274
		1275
		1276
		1277
		1278
		1279
		1280
		1281
		1282
		1283
		1284
		1285
		1286
		1287
		1288
		1289
		1290
		1291
		1292
		1293
		1294
		1295
		1296
		1297
		1298
		1299
		1300
		1301
		1302
		1303
		1304
		1305
		1306
		1307
		1308
		1309
		1310
		1311
		1312
		1313
		1314
		1315
		1316
		1317
		1318
		1319
		1320
		1321
		1322
		1323
		1324
		1325
		1326
		1327
		1328
		1329
		1330
		1331
		1332
		1333
		1334
		1335
		1336
		1337
		1338
		1339
		1340
		1341
		1342
		1343
		1344
		1345
		1346
		1347
		1348
		1349
		1350
		1351
		1352
		1353
		1354
		1355
		1356
		1357
		1358
		1359
		1360
		1361
		1362
		1363
		1364
		1365
		1366
		1367
		1368
		1369
		1370
		1371
		1372
		1373
		1374
		1375
		1376
		1377
		1378
		1379
		1380
		1381
		1382
		1383
		1384
		1385
		1386
		1387
		1388
		1389
		1390
		1391
		1392
		1393
		1394
		1395
		1396
		1397
		1398
		1399
		1400
		1401
		1402
		1403
		1404
		1405
		1406
		1407
		1408
		1409
		1410
		1411
		1412
		1413
		1414
		1415
		1416
		1417
		1418
		1419
		1420
		1421
		1422
		1423
		1424
		1425
		1426
		1427
		1428
		1429
		1430
		1431
		1432
		1433
		1434
		1435
		1436
		1437
		1438
		1439
		1440
		1441
		1442
		1443
		1444
		1445
		1446
		1447
		1448
		1449
		1450
		1451
		1452
		1453
		1454
		1455
		1456
		1457
		1458
		1459
		1460
		1461
		1462
		1463
		1464
		1465
		1466
		1467
		1468
		1469
		1470
		1471
		1472
		1473
		1474
		1475
		1476
		1477
		1478
		1479
		1480
		1481
		1482
		1483
		1484
		1485



1088	Trieu H. Trinh, Yuhuai Wu, Quoc V. Le, He He,	Renqiu Xia, Mingsheng Li, Hancheng Ye, Wenjie Wu,	1146
1089	and Thang Luong. 2024. <a href="#">Solving olympiad ge-</a>	Hongbin Zhou, Jiakang Yuan, Tianshuo Peng, Xinyu	1147
1090	<a href="#">ometry without human demonstrations</a> . <i>Nature</i> ,	Cai, Xiangchao Yan, Bin Wang, Conghui He, Bo-	1148
1091	625(7995):476–482.	tian Shi, Tao Chen, Junchi Yan, and Bo Zhang. 2025.	1149
1092	Aaron van den Oord, Oriol Vinyals, and Koray	<a href="#">Geox: Geometric problem solving through unified</a>	1150
1093	Kavukcuoglu. 2017. Neural discrete representation	<a href="#">formalized vision-language pre-training</a> . In <i>The Thir-</i>	1151
1094	learning. In <i>Proceedings of the 31st International</i>	<i>teenth International Conference on Learning Repre-</i>	1152
1095	<i>Conference on Neural Information Processing Sys-</i>	<i>sentations</i> .	1153
1096	<i>tems</i> , NIPS’17, page 6309–6318, Red Hook, NY,	Liangyu Xu, Yingxiu Zhao, Jingyun Wang, Yingyao	1154
1097	USA. Curran Associates Inc.	Wang, Bu Pi, Chen Wang, Mingliang Zhang, Ji-	1155
1098	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	hao Gu, Xiang Li, Xiaoyong Zhu, Jun Song, and	1156
1099	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	Bo Zheng. 2025. <a href="#">Geosense: Evaluating identification</a>	1157
1100	Kaiser, and Illia Polosukhin. 2017. Attention is all	<a href="#">and application of geometric principles in multimodal</a>	1158
1101	you need. In <i>Proceedings of the 31st International</i>	<a href="#">reasoning</a> . <i>Preprint</i> , arXiv:2504.12597.	1159
1102	<i>Conference on Neural Information Processing Sys-</i>	Shihao Xu, Yiyang Luo, and Wei Shi. 2024. <a href="#">Geo-llava:</a>	1160
1103	<i>tems</i> , NIPS’17, page 6000–6010, Red Hook, NY,	<a href="#">A large multi-modal model for solving geometry</a>	1161
1104	USA. Curran Associates Inc.	<a href="#">math problems with meta in-context learning</a> . In	1162
1105	Vicuna. 2023. Vicuna: An open-source chatbot im-	<i>Proceedings of the 2nd Workshop on Large Genera-</i>	1163
1106	pressing gpt-4 with 90 <a href="https://lmsys.org/blog/2023-03-30-vicuna/">https://lmsys.org/blog/</a>	<i>tive Models Meet Multimodal Applications</i> , LGM3A	1164
1107	<a href="https://lmsys.org/blog/2023-03-30-vicuna/">2023-03-30-vicuna/</a> .	’24, page 11–15, New York, NY, USA. Association	1165
1108	Ke Wang, Juntong Pan, Weikang Shi, Zimu Lu, Houxing	for Computing Machinery.	1166
1109	Ren, Aojun Zhou, Mingjie Zhan, and Hongsheng Li.	Yibo Yan, Jiamin Su, Jianxiang He, Fangteng Fu,	1167
1110	2024a. <a href="#">Measuring multimodal mathematical reason-</a>	Xu Zheng, Yuanhuiyi Lyu, Kun Wang, Shen Wang,	1168
1111	<a href="#">ing with MATH-vision dataset</a> . In <i>The Thirty-eight</i>	Qingsong Wen, and Xuming Hu. 2025. <a href="#">A survey</a>	1169
1112	<i>Conference on Neural Information Processing Sys-</i>	<a href="#">of mathematical reasoning in the era of multimodal</a>	1170
1113	<i>tems Datasets and Benchmarks Track</i> .	<a href="#">large language model: Benchmark, method &amp; chal-</a>	1171
1114	Zhikai Wang, Jiashuo Sun, Wenqi Zhang, Zhiqiang Hu,	<a href="#">lenges</a> . <i>Preprint</i> , arXiv:2412.11936.	1172
1115	Xin Li, Fan Wang, and Deli Zhao. 2025. <a href="#">Benchmark-</a>	Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chun-	1173
1116	<a href="#">ing multimodal mathematical reasoning with explicit</a>	yan Li, and Jianfeng Gao. 2023. <a href="#">Set-of-mark</a>	1174
1117	<a href="#">visual dependency</a> . <i>Preprint</i> , arXiv:2504.18589.	<a href="#">prompting unleashes extraordinary visual grounding</a>	1175
1118	Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen,	<a href="#">in gpt-4v</a> . <i>Preprint</i> , arXiv:2310.11441.	1176
1119	Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Hao-	Weichen Yu, Ziyang Yang, Shanchuan Lin, Qi Zhao,	1177
1120	tian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev	Jiayin Wang, Liangke Gui, Matt Fredrikson, and	1178
1121	Arora, and Danqi Chen. 2024b. <a href="#">Charxiv: Charting</a>	Lu Jiang. 2024. <a href="#">Is your text-to-image model robust</a>	1179
1122	<a href="#">gaps in realistic chart understanding in multimodal</a>	<a href="#">to caption noise?</a> <i>Preprint</i> , arXiv:2412.19531.	1180
1123	<a href="#">LLMs</a> . In <i>The Thirty-eight Conference on Neural</i>	Zhou Yu, Jun Yu, Yuhao Cui, Dacheng Tao, and Qi Tian.	1181
1124	<i>Information Processing Systems Datasets and Bench-</i>	2019. Deep modular co-attention networks for visual	1182
1125	<i>marks Track</i> .	question answering. In <i>Proceedings of the IEEE/CVF</i>	1183
1126	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	<i>Conference on Computer Vision and Pattern Recogn-</i>	1184
1127	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,	<i>ition (CVPR)</i> .	1185
1128	and Denny Zhou. 2022. Chain-of-thought prompt-	Yuan Yuan, Zhaojian Li, and Bin Zhao. 2025. <a href="#">A survey</a>	1186
1129	ing elicits reasoning in large language models. In	<a href="#">of multimodal learning: Methods, applications, and</a>	1187
1130	<i>Proceedings of the 36th International Conference on</i>	<a href="#">future</a> . <i>ACM Comput. Surv.</i> , 57(7).	1188
1131	<i>Neural Information Processing Systems</i> , NIPS ’22,	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng,	1189
1132	Red Hook, NY, USA. Curran Associates Inc.	Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu	1190
1133	Junda Wu, Zhehao Zhang, Yu Xia, Xintong Li,	Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao	1191
1134	Zhaoyang Xia, Aaron Chang, Tong Yu, Sungchul	Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan	1192
1135	Kim, Ryan A. Rossi, Ruiyi Zhang, Subrata Mitra,	Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, and	1193
1136	Dimitris N. Metaxas, Lina Yao, Jingbo Shang, and	3 others. 2024. Mmmu: A massive multi-discipline	1194
1137	Julian McAuley. 2024a. <a href="#">Visual prompting in multi-</a>	<a href="#">multimodal understanding and reasoning benchmark</a>	1195
1138	<a href="#">modal large language models: A survey</a> . <i>Preprint</i> ,	<a href="#">for expert agi</a> . In <i>Proceedings of the IEEE/CVF Con-</i>	1196
1139	arXiv:2409.15310.	<i>ference on Computer Vision and Pattern Recognition</i>	1197
1140	Wenjun Wu, Lingling Zhang, Jun Liu, Xi Tang, Yaxian	<i>(CVPR)</i> , pages 9556–9567.	1198
1141	Wang, Shaowei Wang, and Qianying Wang. 2024b.	Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov,	1199
1142	<a href="#">E-gps: Explainable geometry problem solving via</a>	and Lucas Beyer. 2023. Sigmoid loss for lan-	1200
1143	<a href="#">top-down solver and bottom-up generator</a> . In <i>2024</i>	<a href="#">guage image pre-training</a> . In <i>Proceedings of the</i>	1201
1144	<i>IEEE/CVF Conference on Computer Vision and Pat-</i>	<i>IEEE/CVF International Conference on Computer</i>	1202
1145	<i>tern Recognition (CVPR)</i> , pages 13828–13837.	<i>Vision (ICCV)</i> , pages 11975–11986.	1203

- Jiaxin Zhang and Yashar Moshfeghi. 2024. [GOLD: Geometry problem solver with natural language description](#). In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 263–278, Mexico City, Mexico. Association for Computational Linguistics.
- Ming-Liang Zhang, Fei Yin, Yi-Han Hao, and Cheng-Lin Liu. 2022. [Plane geometry diagram parsing](#). In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 1636–1643. International Joint Conferences on Artificial Intelligence Organization. Main Track.
- Ming-Liang Zhang, Fei yin, and Cheng-Lin Liu. 2023. [A multi-modal neural geometric solver with textual clauses parsed from diagram](#). In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23*, pages 3374–3382. International Joint Conferences on Artificial Intelligence Organization. Main Track.
- Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu Qiao, Peng Gao, and Hongsheng Li. 2025a. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? In *Computer Vision – ECCV 2024*, pages 169–186, Cham. Springer Nature Switzerland.
- Renrui Zhang, Xinyu Wei, Dongzhi Jiang, Ziyu Guo, Yichi Zhang, Chengzhuo Tong, Jiaming Liu, Aojun Zhou, Shanghang Zhang, Peng Gao, and Hongsheng Li. 2025b. MAVIS: Mathematical visual instruction tuning with an automatic data engine. In *The Thirteenth International Conference on Learning Representations*.
- Shan Zhang, Aotian Chen, Yanpeng Sun, Jindong Gu, Yi-Yu Zheng, Piotr Koniusz, Kai Zou, Anton van den Hengel, and Yuan Xue. 2025c. [Open eyes, then reason: Fine-grained visual mathematical understanding in mllms](#). *Preprint*, arXiv:2501.06430.
- Xiaokai Zhang, Na Zhu, Yiming He, Jia Zou, Cheng Qin, Yang Li, and Tuo Leng. 2024a. [Fgeo-sss: A search-based symbolic solver for human-like automated geometric reasoning](#). *Symmetry*, 16(4).
- Xiaokai Zhang, Na Zhu, Cheng Qin, Yang Li, Zhenbing Zeng, and Tuo Leng. 2024b. [Fgeo-hypergnet: Geometric problem solving integrating formal symbolic system and hypergraph neural network](#). *Preprint*, arXiv:2402.11461.
- Xiaokai Zhang, Na Zhu, Cheng Qin, LI Yang, Zhenbing Zeng, and Tuo Leng. 2024c. [Formal representation and solution of plane geometric problems](#). In *The 4th Workshop on Mathematical Reasoning and AI at NeurIPS’24*.
- Xinsong Zhang, Yarong Zeng, Xinting Huang, Hu Hu, Runquan Xie, Han Hu, and Zhanhui Kang. 2025d. [Low-hallucination synthetic captions for large-scale vision-language model pre-training](#). *Preprint*, arXiv:2504.13123.
- Zeren Zhang, Jo-Ku Cheng, Jingyang Deng, Lu Tian, Jinwen Ma, Ziran Qin, Xiaokai Zhang, Na Zhu, and Tuo Leng. 2025e. [Diagram formalization enhanced multi-modal geometry problem solver](#). In *ICASSP 2025 - 2025 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5.
- Junbo Zhao, Ting Zhang, Jiayu Sun, Mi Tian, and Hua Huang. 2025. [Pi-gps: Enhancing geometry problem solving by unleashing the power of diagrammatic information](#). *Preprint*, arXiv:2503.05543.
- Ling Zhong, Yujing Lu, Jing Yang, Weiming Li, Peng Wei, Yongheng Wang, Manni Duan, and Qing Zhang. 2025. [Domaincqa: Crafting expert-level qa from domain-specific charts](#). *Preprint*, arXiv:2503.19498.
- Na Zhu, Xiaokai Zhang, Qike Huang, Fangzhen Zhu, Zhenbing Zeng, and Tuo Leng. 2025. [Fgeo-parser: Autoformalization and solution of plane geometric problems](#). *Symmetry*, 17(1).
- Chengke Zou, Xingang Guo, Rui Yang, Junyu Zhang, Bin Hu, and Huan Zhang. 2025. [Dynamath: A dynamic visual benchmark for evaluating mathematical reasoning robustness of vision language models](#). In *The Thirteenth International Conference on Learning Representations*.
- Jia Zou, Xiaokai Zhang, Yiming He, Na Zhu, and Tuo Leng. 2024. [Fgeo-drl: Deductive reasoning for geometric problems through deep reinforcement learning](#). *Symmetry*, 16(4).

## A Additional axis on benchmark dataset

### A.1 Reasoning complexity

We discuss the mathematical concepts and difficulty levels encountered in plane geometry problems used by existing benchmarks and datasets. Typical plane geometry problems involve calculating specific angle measures, arc measures, segment or arc lengths, and areas of designated regions. Computing these numerical values generally requires basic arithmetic and root operations, but may also involve trigonometric functions, such as sine and cosine. Although no standardized quantitative method currently exists to measure problem difficulty, problems can be qualitatively categorized according to their original sources, such as SAT exams (Seo et al., 2015; Sachan et al., 2017; Sachan and Xing, 2017), plane geometry curricula from grades 6–12 American (Lu et al., 2021; Zhang et al., 2023; Sun et al., 2024) or Chinese school (Chen et al., 2021; Cao and Xiao, 2022; Xu et al., 2025), college-level mathematics (Yue et al., 2024), or mathematics competitions, e.g., AMC 8, 10, and 12 (Wang et al., 2024a).

### A.2 Diagram-text redundancy

To serve as rigorous benchmarks and datasets for multi-modal reasoning, the collected problems must require simultaneous interpretation of both diagrams and accompanying textual descriptions. By contrast, PGPS problems that can be solved using the text alone cannot effectively evaluate the diagram-text integration capability of PGPS methods. Nevertheless, many existing benchmarks and datasets still contain such problems, thereby inadequately assessing the perception abilities of PGPS methods (Zhang et al., 2025a).

Recent PGPS benchmarks have addressed this limitation by explicitly annotating problems with modality-specific information and subsequently removing redundant textual cues (Lu et al., 2021; Zhang et al., 2023, 2025a). Several benchmarks provide multiple variants of each problem for more fine-grained analysis of diagram-text dependency. For instance, MathVerse (Zhang et al., 2025a) relocates selected information from the text into the diagram, while DynaMath (Zou et al., 2025) generates alternative diagrams and corresponding answers based on a single textual description. Thus, failure to solve certain variants of the same problem indicates that the model is not genuinely utilizing the diagram.

### A.3 Data collection methods

We summarize three data collection methods mainly used to construct PGPS datasets.

**Human annotation** In most cases, datasets are constructed through human annotation based on problems sourced from textbooks, internet sites, or similar resources (Seo et al., 2015; Chen et al., 2021; Lu et al., 2021, 2024; Sun et al., 2024; Yue et al., 2024). This involves manually collecting problems and having human annotators provide the corresponding outputs. Additionally, some studies apply text augmentation techniques, such as back-translation, to diversify the text style and enrich the dataset (Cao and Xiao, 2022).

**Synthetic annotation** Several PGPS studies create synthetic benchmarks and datasets instead of collecting problems from textbooks or the internet. These studies typically implement synthetic engines to generate diagrams and corresponding structured information. For example, synthetic engines can generate captions containing the geometric information explicitly present in diagrams (Zhang et al., 2025b), or use symbolic reasoning engines to produce reasoning steps that derive the stated goals from diagram-text pairs (Zhang et al., 2025b; Kazemi et al., 2024; Fu et al., 2025). Such synthetic approaches offer clear advantages, including easy scalability and guaranteed completeness of annotations. However, they often struggle to produce sufficiently diverse diagrams that accurately reflect the real-world problems. This limitation is further discussed in §6.2.

**L(V)LM-assisted annotation** For certain datasets, particularly those with natural-language description as the output representation, LLMs and VLMs such as GPT (Brown et al., 2020) or GPT-4V (OpenAI, 2023) are employed for dataset construction. Specifically, problems and solutions are sourced from datasets like GeoQA+, UniGeo, or PGPS9K, and GPT or GPT-4V are used to augment these by generating multiple problem-solution pairs for a given problem scenario (Gao et al., 2025; Shi et al., 2024; Zhang et al., 2025b). Alternatively, some studies apply the same process to synthetic data, such as diagram-caption pairs generated by a synthetic data engine (Zhang et al., 2025b; Kazemi et al., 2024). However, due to the poor perception ability of GPT-4V, several hallucinations occur in the augmented datasets. We discuss more details about the challenge in §6.1.

## B PGPS Methods

### B.1 Summary of PGPS methods

We summarize the PGPS methods in terms of the encoder, intermediate representation, decoder, and the output format at Table A1.

## C Challenges and Future Directions

### C.1 Error analysis on wrong responses

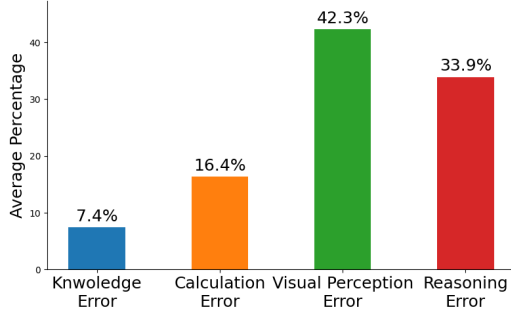


Figure A1: Error analysis on the response of GPT-4V on MathVerse. We analyze the responses of GPT-4V on MathVerse, reporting the average percentage for each type of error across five MathVerse variants, Text Dominant, Text Lite, Vision Intensive, Vision Dominant, and Vision Only, which are reported in MathVerse. Our analysis indicates that incorrect answers predominantly result from visual perception and reasoning errors.

### C.2 Examples of perception hallucinations

We provide examples of hallucinated responses by GPT-4.1 in Table A2.

### C.3 Comprehensivity of current PGPS benchmarks

Methods	Realistic styles of diagrams	No data leakage	Diagram-text interdependence
MMMU	○	○	×
Math-V	○	○	×
MathVista	○	×	×
MathVerse	○	×	○
GeomVerse	×	○	×
VisOnlyQA	×	○	○
MM-Math	○	○	×
GeoEval	×	×	×
DynaMath	○	×	○

Table A3: Comprehensivity across existing PGPS benchmarks. The table summarizes benchmark features in terms of realistic diagram styles, absence of data leakage, and consideration of diagram-text interdependence.

### C.4 Synthetic and real-world geometric diagrams

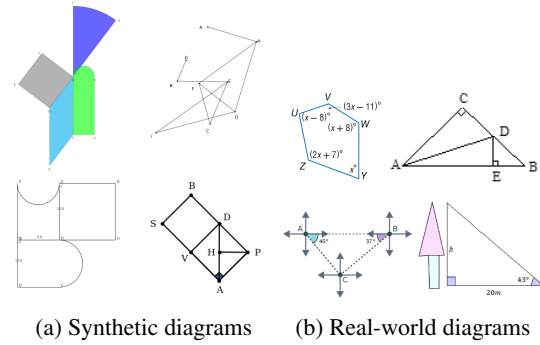


Figure A2: Visualization of the synthetic and real-world geometric diagrams. We compare the geometric diagrams, which are synthetically generated or manually collected from existing sources. The synthetic diagrams are from GeomVerse, VisOnlyQA, MAVIS, and GeoDANO. The real-world diagrams are from MathVerse.



Encoder	Intermediate	Decoder	Output	Methods
E1	I1	–	–	HoughGeo (Chen et al., 2015), G-Aligner (Seo et al., 2014), GEOS (Seo et al., 2015)
E2	I1	–	–	PGDPNet (Zhang et al., 2022), FGeo-Parser (Zhu et al., 2025)
E1	I1	D1	O1	GEOS++ (Sachan et al., 2017), GEOS-OS (Sachan and Xing, 2017), GeoShader (Alvin et al., 2017), S2 (Gan et al., 2019)
E2	I1	D2	O1	FGeo-HyperGNet (Zhang et al., 2024b), GCN-GPS (Jian et al., 2023a), GeoDRL (Peng et al., 2023)
E2	I1	D3	O1	InterGPS (Lu et al., 2021), E-GPS (Wu et al., 2024b), Pi-GPS (Zhao et al., 2025), FGeo-DRL (Zou et al., 2024)
E2	I1	D1	O1	FGeo-SSS (Zhang et al., 2024a)
E2	I1	D2	O3	GOLD (Zhang and Moshfeghi, 2024)
E2	I2	D3	O2	NGS (Chen et al., 2021), DPE-NGS (Cao and Xiao, 2022), Geoformer (Chen et al., 2022), PGPSNet (Zhang et al., 2023), SCA-GPS (Ning et al., 2023), UniMath (Liang et al., 2023), FLCL-GPS (Jian et al., 2023b), LANS (Li et al., 2024), GeoX (Xia et al., 2025), GeoDANO (Cho et al., 2025)
E2	I2	D3	O3	Math-LLaVA (Shi et al., 2024), Visual SKETCH-PAD (Hu et al., 2024b), MAVIS (Zhang et al., 2025b), G-LLaVA (Gao et al., 2025), DFE-GPS (Zhang et al., 2025e), Chimera (Peng et al., 2025), Geo-LLaVA (Xu et al., 2024), SVE-Math (Zhang et al., 2025c)

Table A1: Categorization of existing PGPS methods. We categorize the PGPS methods based on their encoder, intermediate representation, decoder, and output format. The symbols come from Fig. 3.

	Example 1	Example 2
<b>Input prompt</b>	In the given diagram, extend the segment SU to form an infinite straight line. This line divides the two-dimensional plane into two half-planes. The points Q and A are also shown in the diagram. Are the points Q and A located on the same side of the line SU?	
<b>Input diagram</b>		
<b>Response</b>	<p>To determine if points Q and A are on the same side of the line SU, observe the diagram: extending SU clearly shows that point Q is on one side of the line, <b>while point A is on the opposite side.</b> The quadrilateral QASU suggests <b>A and Q are separated by the segment SU itself.</b></p>	<p>To determine whether points Q and A are on the same side of the line SU, we observe the diagram. The segment SU is extended to form a straight line. Point Q is above SU, <b>while point A is on the opposite side</b>, to the left of SU. Therefore, <b>the points Q and A are on opposite sides of line SU.</b></p>

Table A2: Examples of hallucinations occurred in GPT-4.1 responses. We analyzed the diagram perception performance of GPT-4.1 (OpenAI, 2025a), specifically determining whether two points are on the same side of a given line. We generated 100 problems using the synthetic data engine from GeoDANO (Cho et al., 2025) and tested them with GPT-4.1, observing a low accuracy of 59%. The examples above illustrate cases where GPT-4.1 hallucinated and provided incorrect answers. Hallucinated parts are highlighted in bold.