# VisOnlyQA: Large Vision Language Models Still Struggle with Visual Perception of Geometric Information

Ryo Kamoi, Yusen Zhang, Sarkar Snigdha Sarathi Das, Ranran Haoran Zhang, Rui Zhang

{ryokamoi, rmz5227}@psu.edu          Penn State University          https://visonlyqa.github.io/

## Motivation

The ability of LVLMs to perceive visual information in images has not been sufficiently studied. Specifically, it remains unclear **how accurately LVLMs can perceive geometric information, such as shape, angle, and size**, while geometric perception is fundamental to understanding visual information.

Limitations in existing datasets:

1. Popular datasets, such as MMMU and MathVista, target tasks that require expert-level reasoning and knowledge. They are not suitable for analyzing perception.

2. Datasets for evaluating LVLMs at perceiving visual information include high-level tasks, such as scene understanding, which do not necessarily require accurate geometric perception.

## Dataset Creation

We propose *VisOnlyQA*, a new dataset designed to evaluate how accurately LVLMs can perceive basic geometric information in images. Our dataset includes questions that **directly ask about basic and common geometric information (e.g., length, angle, and shape)**.

VisOnlyQA is designed to have favorable properties for analyzing the capability of LVLMs to perceive geometric information:
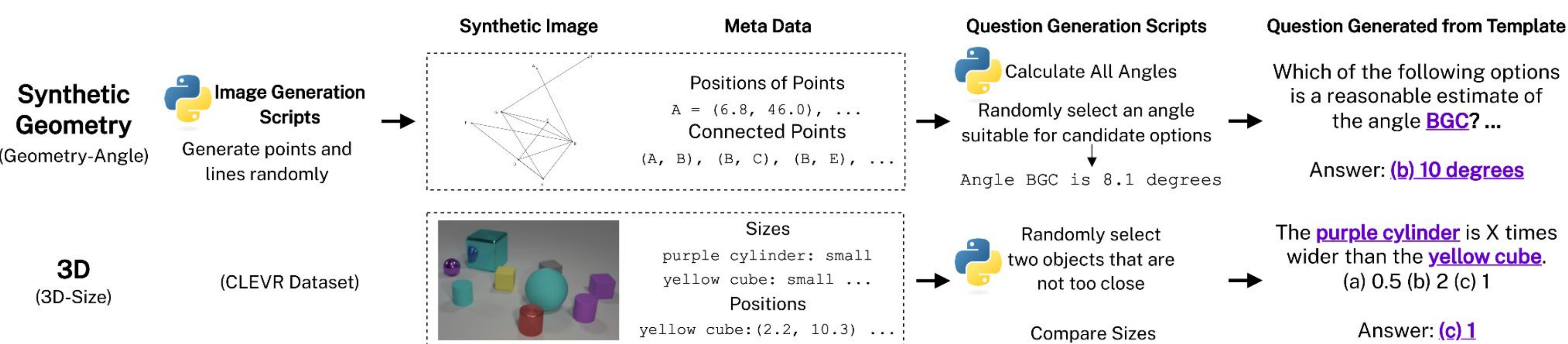
1. The questions in our dataset do not involve challenging reasoning or knowledge, exclusively evaluating geometric perception.

2. We use scientific figures to create unambiguous questions that require accurate geometric perception.

### Real and Synthetic Figures
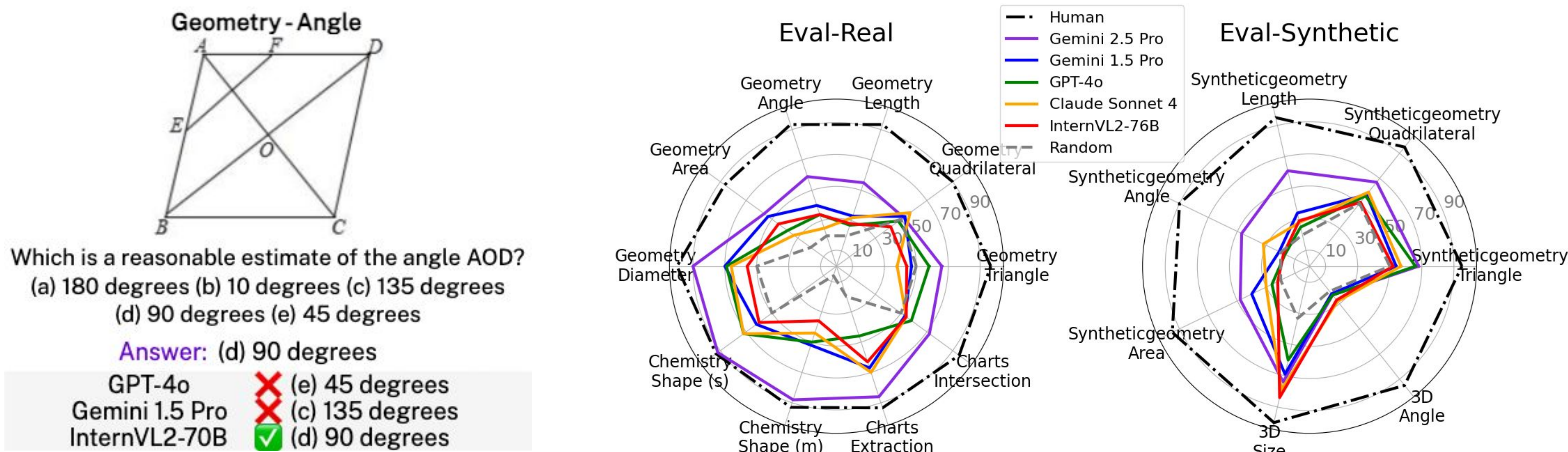
VisOnlyQA consists of two types of figures.

**Real Figures**: Images in popular datasets, such as MathVista and MMMU. We manually annotate new questions.

**Synthetic Figures**: We generate synthetic figures where we can get accurate geometric information like positions and size.



---

**TL;DR** - We introduce *VisOnlyQA*, a dataset for evaluating the geometric perception of LVLMs.

**We observe that even state-of-the-art LVLMs including GPT-4o and Gemini 2.5 Pro struggle with basic geometric perception questions like *"Does this figure include triangle ABC?"***



Geometry - Angle

Which is a reasonable estimate of the angle AOD?
(a) 180 degrees (b) 10 degrees (c) 135 degrees
(d) 90 degrees (e) 45 degrees

Answer: (d) 90 degrees

| | | |
|---|---|---|
| GPT-4o | ✗ | (e) 45 degrees |
| Gemini 1.5 Pro | ✗ | (c) 135 degrees |
| InternVL2-70B | ✓ | (d) 90 degrees |

**VisOnlyQA consists of simple visual perception tasks, on which human performance is nearly perfect. However, state-of-the-art LVLMs such as GPT-4o, Claude Sonnet 4, and Gemini 2.5 Pro perform poorly.**
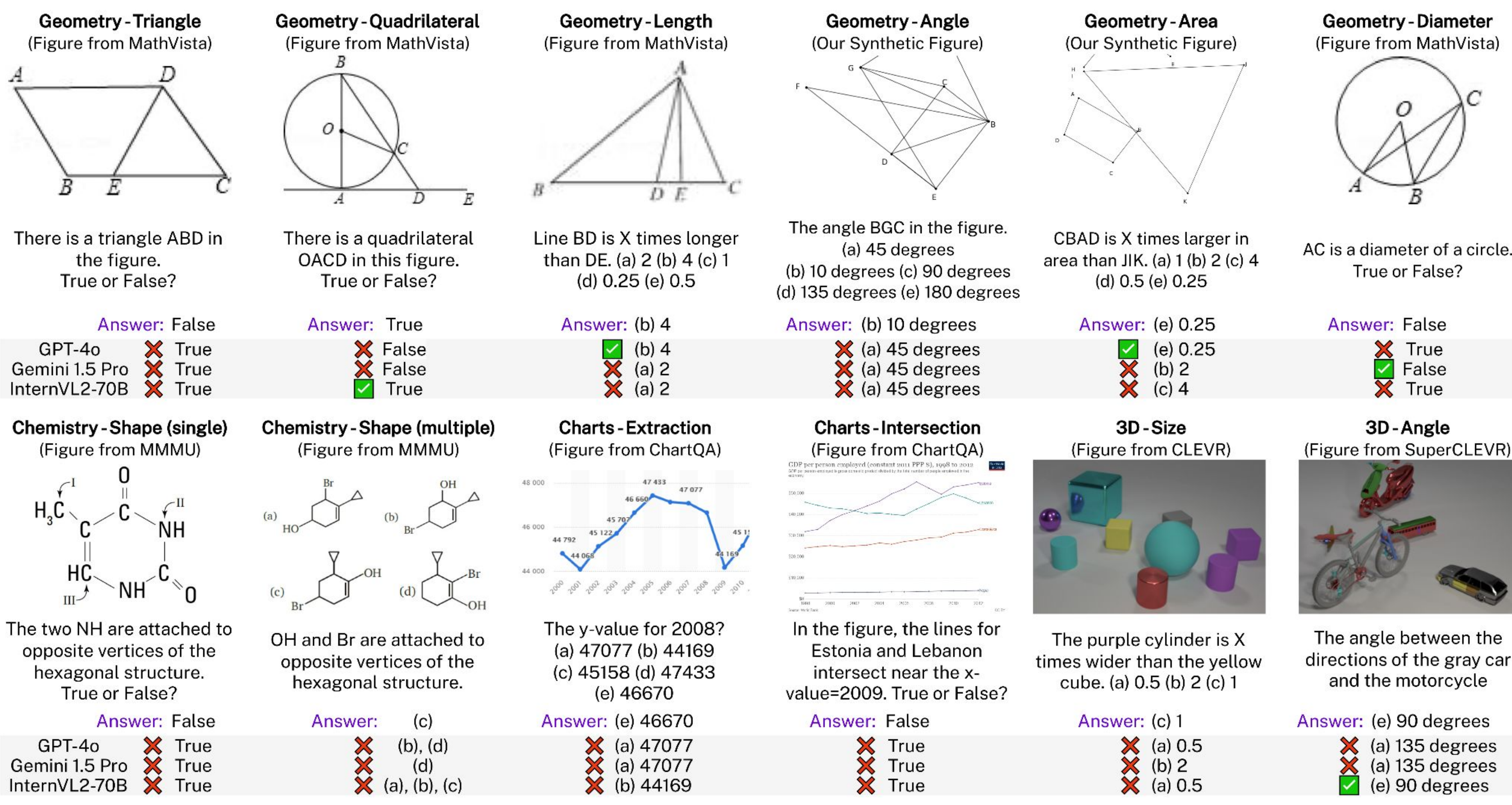
## VisOnlyQA Dataset

VisOnlyQA includes **12 tasks** in total

| | |
|---|---|
| Geometric Shapes | 5 Tasks |
| Chemical Structure | 2 Tasks |
| Charts | 2 Tasks |
| 3D shapes | 2 Tasks |

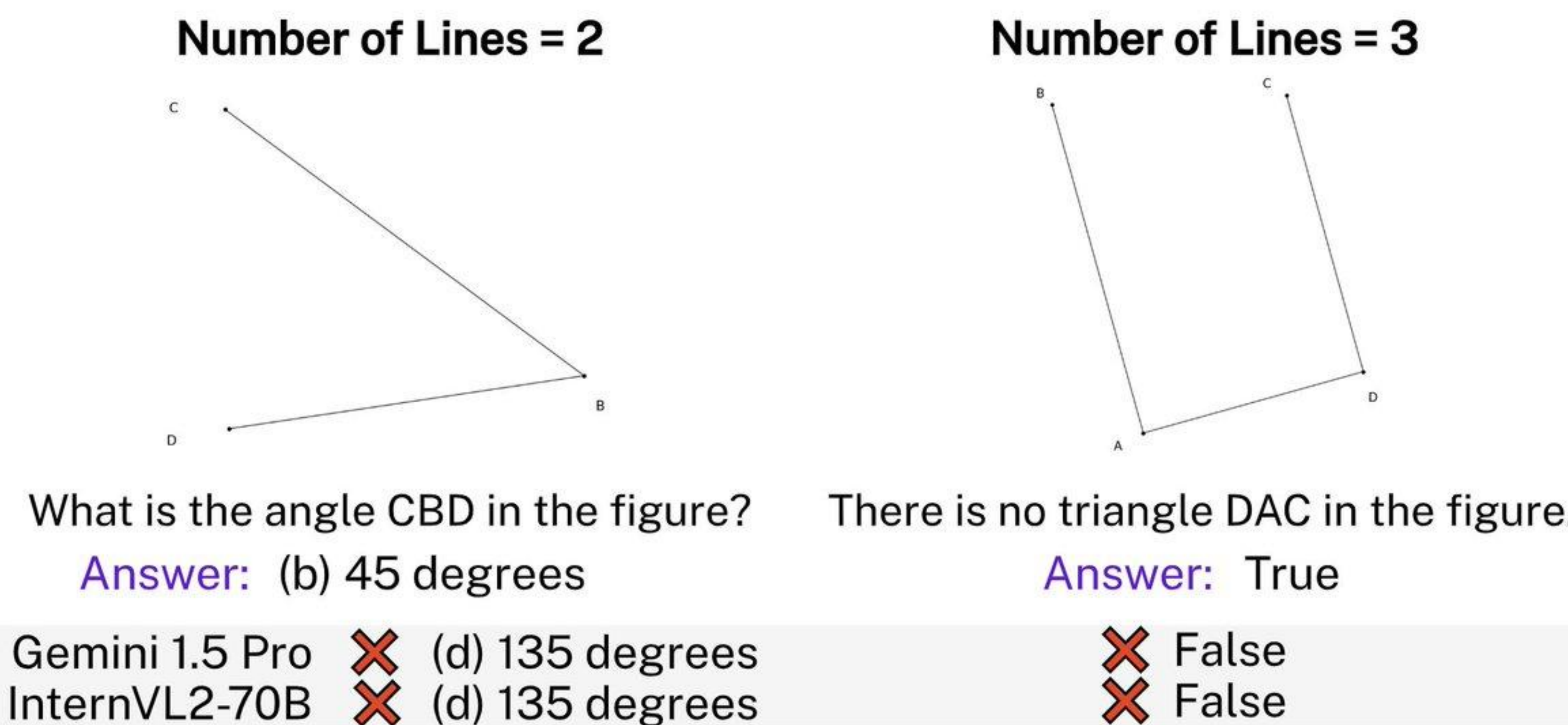*VisOnlyQA* includes the evaluation and training set, consisting of the Real and Synthetic figures

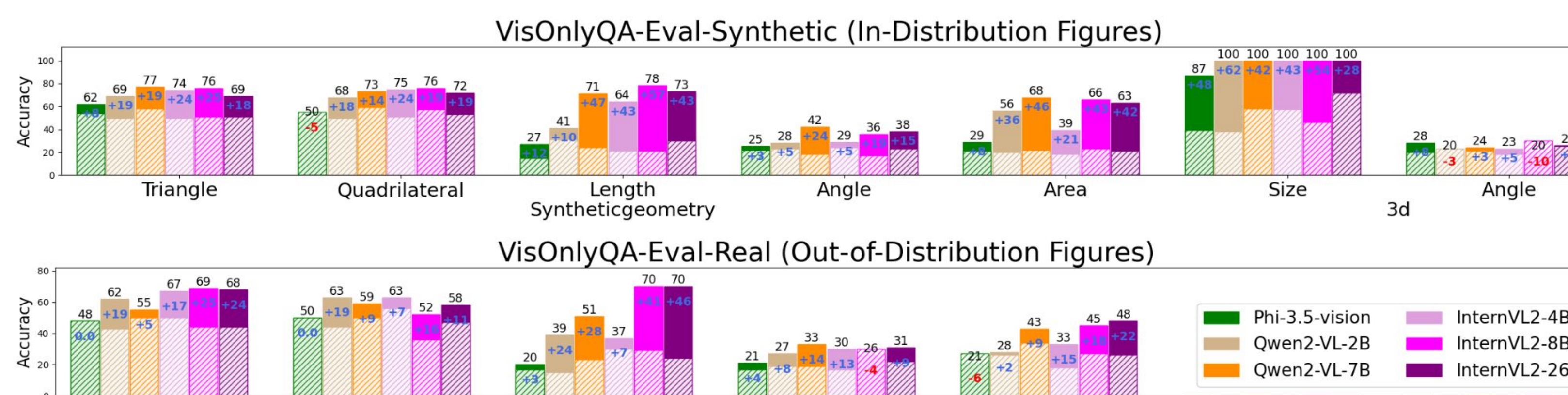| | # Tasks | # Questions |
|---|---|---|
| Eval-Real | 10 Tasks | 900 |
| Eval-Synthetic | 7 Tasks | 700 |
| Train | 7 tasks | 70k |



## Analysis

### Geometric Shapes with Only Two Lines are Already Difficult

Surprisingly, **LVLMS perform poorly even on very simple geometric shapes that only include two or three lines**.



### Fine-tuning Does Not Fully Solve This Issue

**Fine-tuning on VisOnlyQA-Train does not always improve geometric perception even on in-distribution images**. This result suggests that simply scaling training data does not solve this issue



### Larger LMs Improve the Geometric Perception of LVLMs

LVLMs with **larger language models exhibit better visual perception**. This result suggests that language models are crucial in processing visual information encoded by ViT.

| | ViT Size | LLM Size | Real | Synthetic |
|---|---|---|---|---|
| InternVL2-4B | 304M | 3.8B | 38.4 | 34.1 |
| InternVL2-8B | 304M | 7.7B | **40.7** | **35.0** |
| Qwen2-VL-2B | 675M | 1.5B | 32.3 | 33.6 |
| Qwen2-VL-7B | 675M | 7.6B | **38.9** | **37.1** |