# Multimodal Self-Instruct: Synthetic Abstract Image and Visual Reasoning Instruction Using Language Model

**Anonymous ACL submission**

## Abstract

Although most current large multimodal models (LMMs) can already understand photos of natural scenes and portraits, their understanding of abstract images, e.g., charts, maps, or layouts, and visual reasoning capabilities remains quite rudimentary. They often struggle with simple daily tasks, such as reading time from a clock, understanding a flowchart, or planning a route using a road map. In light of this, we design a multi-modal self-instruct, utilizing large language models and their code capabilities to synthesize massive abstract images and visual reasoning instructions across daily scenarios. Our strategy effortlessly creates a multimodal benchmark with 11,193 instructions for eight visual scenarios: charts, tables, simulated maps, dashboards, flowcharts, relation graphs, floor plans, and visual puzzles. **This benchmark, constructed with simple lines and geometric elements, exposes the shortcomings of most advanced LMMs** like GPT-4V and Llava in abstract image understanding, spatial relations reasoning, and visual element induction. Besides, to verify the quality of our synthetic data, we fine-tune an LMM using 62,476 synthetic chart, table and road map instructions. The results demonstrate improved chart understanding and map navigation performance, and also demonstrate potential benefits for other visual reasoning tasks. Our code and data are available at this anonymous link: https://anonymous.4open.science/r/self-instruct-data-engine-E785/

## 1 Introduction

In recent times, spurred by breakthroughs in large language models (LLMs) (Zeng et al., 2023; Touvron et al., 2023a; OpenAI, 2022, 2023; Touvron et al., 2023b; Bi et al., 2024; Jiang et al., 2024; Anthropic, 2024; Abdin et al., 2024), large multimodal models (LMMs) have also undergone rapid advancements (Liu et al., 2024b,a; Team et al., 2023; Bai et al., 2023a; Lu et al., 2024; McKinzie et al., 2024). Leveraging a pre-trained LLM to unify the encoding of all modalities empowers LMMs to understand human daily environments and execute complex tasks (Hong et al., 2023; Hu et al., 2023; Zhang et al., 2023; Koh et al., 2024; Zhang et al., 2024b). This greatly expands the potential of general-purpose AI assistants.

Despite these achievements, LMMs still exhibit significant deficiencies when deployed in human daily life (Yin et al., 2023; Xie et al., 2024). For instance, LMMs often fail when planning a route using a road map, reading the time from a clock image, or interpreting a flowchart. We observe that these simple daily activities require LMMs to understand abstract images, such as maps, charts, and dashboards, rather than natural photographs or portraits with explicit semantics. These abstract images composed of simple geometric elements are more challenging for LMMs. Furthermore, even many advanced LMMs are easily stumped by simple visual-level reasoning tasks, such as geometric pattern induction and visual symbol comparison.

However, these capabilities, i.e., perceiving abstract images and reasoning about visual elements, are essential for LMMs if we deploy an LMM-driven agent in our daily lives. It can help us with data analysis, map navigation, web searches, and many other tedious tasks. On the one hand, despite valuable explorations by some pioneers (Yu et al., 2023b; Liu et al., 2023; Han et al., 2023; Ying et al., 2024; Wei et al., 2024), these abstract image understanding and visual reasoning abilities have not been adequately emphasized, and we need a dedicated benchmark to systematically evaluate the performance of current LMMs in this aspect. On the other hand, unlike semantic-related tasks, collecting such abstract image-text pairs with reasoning context is labor-intensive and time-consuming.

To fill in the gap, we drew inspiration from synthetic data (Wang et al., 2022b; Liu et al., 2024c;
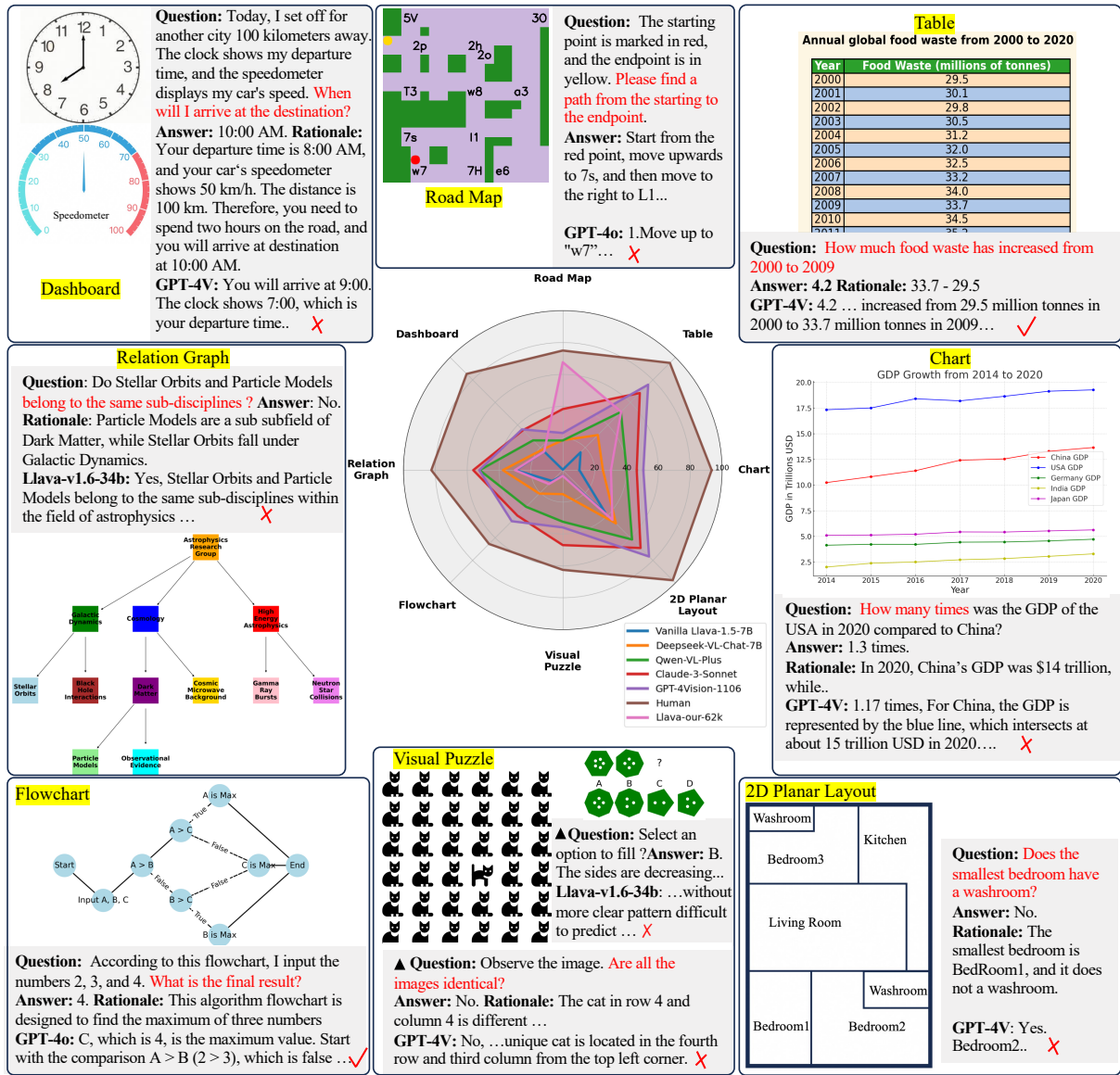
**Dashboard**

**Question:** Today, I set off for another city 100 kilometers away. The clock shows my departure time, and the speedometer displays my car's speed. When will I arrive at the destination?
**Answer:** 10:00 AM. **Rationale:** Your departure time is 8:00 AM, and your car's speedometer shows 50 km/h. The distance is 100 km. Therefore, you need to spend two hours on the road, and you will arrive at destination at 10:00 AM.
**GPT-4V:** You will arrive at 9:00. The clock shows 7:00, which is your departure time.. ✗

**Road Map**

**Question:** The starting point is marked in red, and the endpoint is in yellow. Please find a path from the starting to the endpoint.
**Answer:** Start from the red point, move upwards to 7s, and then move to the right to L1...
**GPT-4o:** 1.Move up to "w7"… ✗

**Table**
Annual global food waste from 2000 to 2020

| Year | Food Waste (millions of tonnes) |
|------|---------------------------------|
| 2000 | 29.5 |
| 2001 | 30.1 |
| 2002 | 29.8 |
| 2003 | 30.5 |
| 2004 | 31.2 |
| 2005 | 32.0 |
| 2006 | 32.5 |
| 2007 | 33.2 |
| 2008 | 34.0 |
| 2009 | 33.7 |
| 2010 | 34.5 |
| 2011 | 35.2 |

**Question:** How much food waste has increased from 2000 to 2009
**Answer:** 4.2 **Rationale:** 33.7 - 29.5
**GPT-4V:** 4.2 … increased from 29.5 million tonnes in 2000 to 33.7 million tonnes in 2009… ✓

**Chart**
GDP Growth from 2014 to 2020
(China GDP, USA GDP, Germany GDP, India GDP, Japan GDP)

**Question:** How many times was the GDP of the USA in 2020 compared to China?
**Answer:** 1.3 times.
**Rationale:** In 2020, China's GDP was $14 trillion, while..
**GPT-4V:** 1.17 times, For China, the GDP is represented by the blue line, which intersects at about 15 trillion USD in 2020…. ✗

**Relation Graph**

**Question:** Do Stellar Orbits and Particle Models belong to the same sub-disciplines ? **Answer:** No.
**Rationale:** Particle Models are a sub subfield of Dark Matter, while Stellar Orbits fall under Galactic Dynamics.
**Llava-v1.6-34b:** Yes, Stellar Orbits and Particle Models belong to the same sub-disciplines within the field of astrophysics … ✗

**Flowchart**

**Question:** According to this flowchart, I input the numbers 2, 3, and 4. What is the final result?
**Answer:** 4. **Rationale:** This algorithm flowchart is designed to find the maximum of three numbers
**GPT-4o:** C, which is 4, is the maximum value. Start with the comparison A > B (2 > 3), which is false … ✓

**Visual Puzzle**

▲**Question:** Select an option to fill ?**Answer:** B. The sides are decreasing...
**Llava-v1.6-34b:** …without more clear pattern difficult to predict … ✗
▲ **Question:** Observe the image. Are all the images identical?
**Answer:** No. **Rationale:** The cat in row 4 and column 4 is different …
**GPT-4V:** No, …unique cat is located in the fourth row and third column from the top left corner. ✗

**2D Planar Layout**
(Washroom, Kitchen, Bedroom3, Living Room, Bedroom1, Washroom, Bedroom2)

**Question:** Does the smallest bedroom have a washroom?
**Answer:** No.
**Rationale:** The smallest bedroom is BedRoom1, and it does not a washroom.
**GPT-4V:** Yes. Bedroom2.. ✗

Legend: Vanilla Llava-1.5-7B; Deepseek-VL-Chat-7B; Qwen-VL-Plus; Claude-3-Sonnet; GPT-4Vision-1106; Human; Llava-our-62k

Figure 1: We leverage LLM and code to synthesize abstract images and self-instruct diverse reasoning instructions, e.g., charts, road maps, dashboards, visual puzzles, and relation graphs. Unlike natural landscapes and human photos, these non-natural images constructed with geometric elements require stronger perception and spatial relation reasoning. Our benchmark indicates that current LMMs are far from human-level performance. They even fail to complete simple daily tasks, e.g., reading the time on a clock or planning a route using a map.

Han et al., 2023; Du et al., 2023), which is widely used to supplement the insufficiency of instruction-following data. For instance, distilling high-quality dialogue data from a strong LLM (Wang et al., 2022b; Xu et al., 2023a; Yu et al., 2023a; Chen et al., 2023a; Zhao et al., 2023), or using external tools to refine the quality of synthetic data (Wei et al., 2023; Lee et al., 2024). However, synthesizing image-text data for LMM is not easy, as current LLMs can not directly generate images. An intuitive approach is to combine LLMs with a text-to-image model for producing <image, question, answer> (Li et al., 2023c; Wu et al., 2023b), but most text-to-image models fail to finely control the details of the image, potentially leading to a misalignment between image and text.

Considering that abstract images are composed of lines and geometric elements, we can utilize code to accurately synthesize them. In light of this, we advocate a code-centric self-instruct strategy to synthesize massive abstract images with reasoning questions and answer pairs. We first instruct LLM to autonomously propose a creative visual idea for a daily scenario and then self-propose the necessary data and code to draw an abstract image, such as plotting a relation graph or house layout. Af-

ter synthesizing images, our strategy self-instructs multiple reasoning question-answer pairs based on the plotting idea and code. This code-centric design can effortlessly synthesize diverse abstract images and reasoning instructions, involving chart interpretation, spatial relation reasoning, visual puzzles, and mathematical geometry problems, and also provide accurate answers and rationale.

As shown in Figure 1, our strategy synthesized an abstract image benchmark for daily scenarios, including 11,193 high-quality instructions covering eight scenarios: Dashboard, Road Map, Chart, Table, Flowchart, Relation Graph, Visual Puzzles, and 2D Planar Layout. Empowered by this benchmark, we evaluate several representative LMMs and identify their significant deficiencies in abstract image understanding and visual reasoning. For example, in the dashboard scene, the best-performing LMM (GPT-4V) only achieved a score of 36.2, far below the human level of 85.3. Besides, to verify the quality of the synthesized data, we also synthesized 62,476 charts and road map instructions for fine-tuning Llava-1.5-7B. Experimental results show that our synthesized data can significantly enhance in-domain performance and also benefit other abstract image reasoning tasks.

Our contributions can be summarized as follows:

- We identify that current LMMs have a significant gap compared to humans in understanding and visually reasoning about abstract images, such as maps, charts, and layouts.

- Utilizing LLM and code, We design a multimodal self-instruct strategy to synthesize a diverse set of abstract images and reasoning instructions, providing value data for LMMs.

- We synthesized a benchmark of 11,193 high-quality abstract images, covering eight common scenarios. Our benchmark reveals significant deficiencies even in advanced LMMs. Besides, we synthesized 62,476 chart and road map instructions for fine-tuning, verifying the effectiveness of the synthesized data.

## 2 Multi-modal Self-Instruct

### 2.1 Overview

Our multi-modal self-instruct is an LLM-driven data synthesis strategy capable of producing abstract images and aligned reasoning instructions for various daily scenarios, including road maps, dashboards, 2D planar layouts, charts, relation graphs, flowcharts, and visual puzzles.

Firstly, our strategy can autonomously propose a creative idea for visual scenarios, e.g., *using a step-by-step flowchart to demonstrate how to attend an academy conference* or *designing road map* (Section 2.2). Then it generates detailed code to visualize this idea (Section 2.3). After synthesizing the desired image, LLMs self-instruct multiple high-quality Q&A pairs for this visual content (Section 2.4). The entire process is fully completed by the LLM with a few demonstrations.

As shown in Figure 2, we illustrate the entire process of our image-text synthesis, including using road maps for navigation, interpreting pie charts, solving visual puzzles, and using operating workflow. For each scenario, we synthesize multiple questions, annotated answers, and rationales. For example, in the pie chart case, the LLM designs a multi-step math question about the difference between the largest and smallest categories.

### 2.2 Visual Idea Proposal

To generate an image from scratch, we first instruct the LLM to propose an innovative visual idea. This visual idea illustrates a scenario commonly encountered in daily life or work, e.g., a chart about a specific topic or a road map. Besides, this scenario image can be rendered with code, rather than real portraits or natural scenes. Therefore, we focus on eight common types of abstract images that are rarely covered in current datasets:

```
Working Scene and Life Scene
Charts and Table: Line, bar, pie, composite
 charts, and single and multiple tables.
Flowchart: Algorithm flowcharts and
operating workflows, such as designing a
 slide presentation.
Relation Graph: Multiple relational graphs
 with complex connections.
Road Map: Simulated road maps annotated
with intersection names.
Visual Puzzles: 1. Inductive reasoning
across multiple images. 2. Comparing the
 differences between multiple images.
2D Planar Layout: Floor plans with
different structures and layouts.
Instrument Dashboards: Mechanical dials,
such as clocks, odometers, speedometers,
 thermometers, barometers..
```

We design some examples for each scenario as in-context demonstrations. Prompted by them, the LLM is encouraged to propose a creative and detailed plotting idea using natural language. These visual ideas depict the basic outlines of visual in-
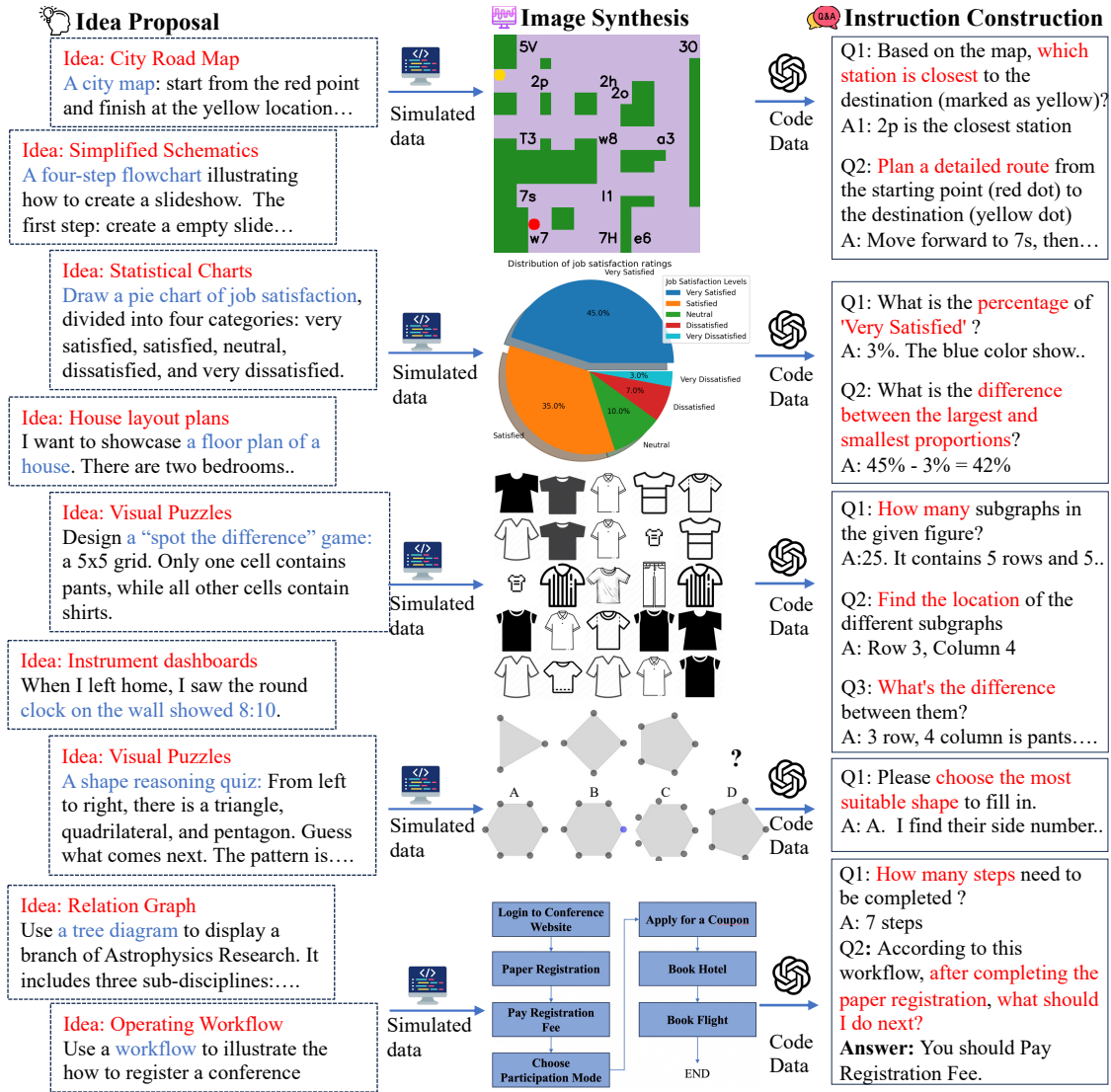
Figure 2: Our multi-modal self-instruct strategy first self-proposes a visual idea to depict an abstract image. Based on this, the LLM generates simulated data and writes code to create the drawings. Subsequently, LLM is instructed to design multiple Q&A based on the code and idea, covering various aspects such as spatial reasoning, color recognition, and mathematical reasoning, constructing a rich set of multimodal instructions.

formation. By incorporating detailed parameters, a visual idea can control the specifics of image synthesis, enabling the creation of a diverse range of images. Additionally, when constructing visual instructions, visual ideas can provide a visual reference for the generation of instructions in natural language form.

## 2.3 Image Synthesis

**Simulated Data** To render the proposed idea into an image, we guide the LLM to first generate some simulated data for the proposed idea. For example, for the pie chart in Figure 2, the LLM needs to fabricate the percentage data for the four types.

**Code Generation** After producing simulated data, LLM generates corresponding Python code to visualize the proposed idea. We encourage the LLM to use popular visualization packages, e.g., Matplotlib[1] or ECharts[2], to create desired visual elements, as it significantly reduces the complexity of code generation. Besides, we instruct the LLM to explicitly define all parameters in the code for plotting images, such as image style, color, font size, and legend position. These explicitly stated parameters control the details of the synthesized images and can be used to produce Q&A.

---

[1]https://matplotlib.org
[2]https://echarts.apache.org/zh/index.html

## 2.4 Visual Instruction Construction

After executing the code, we obtain the expected image. Next, the LLM autonomously proposes multiple high-quality <question, answer> pairs related to this synthetic image.

**Question-Answer Pair Generation.** To make the LLM aware of all the image details, we concatenate the proposed idea, simulated data, and generated code in the prompt, and then guide the LLM to design instructions following data for this synthesized image. More than just image comprehension and captioning tasks, our strategy can self-propose a wide range of unconventional questions for this synthesized image, such as comparing differences among multiple images, area estimation, and spatial relation inference. Furthermore, it can even design diverse multi-step reasoning problems based on multiple synthesized images.

**Annotate Answers with Rationale.** To enhance the training effectiveness of multimodal instruction-following data, we also provide a detailed rationale for each question. We prompt the LLM to carefully review the idea and code, and then generate a detailed rationale for the given question, rather than just providing an answer. Similar to the chain-of-thought process, rationale can be used to train LMMs, enhancing their reasoning capabilities.

Below is a complete case for our pipeline, including Idea Proposal, Image Synthesis, and Instruction Construction. We also provide the results of GPT-4 and Gemini-1.5, which all failed on this case.

```
Idea Proposal: Draw a clock with hour and
minute hands.
Simulated Data: time='8:10', Shape='Round
Clock', color='black', size=...
Code Generation: 'import pyechart...'
Instruction Construction
Question: What time is shown on the dial?
Answer1: 8:10
GPT-4V: 10:10. Gemini-1.5-pro: 2:42.
Math Question: When I left home, the clock
 showed the time indicated in the figure
. What time is it after 8 hours of work?
Answer2: 4:10 or 16:10
Rationale: I see that the clock shows the
time as 8:10. After working for eight
hours, the time should be 16:10.
GPT-4V: 7:10. The clock shows 11:10 ...
Gemini-1.5-pro: 9:50. The time is 1:50 ...
Reasoning Question: I exercised for one and
 a half hours. After finishing, the
clock showed the time as illustrated.
What number did the hour hand point to
when I started my workout?
Answer3: 6 or 7
Rationale: I read the time from the clock
as 8:10, and you have been exercising
```

| Task | #Image | # Instruction | #Usage |
|------|--------|---------------|--------|
| Chart | 1,768 | 34,590 | Train |
| Table | 570 | 10,886 | Train |
| Road map | 17,000 | 17,000 | Train |
| **All** | 19,338 | 62,476 | Train |
| Chart | 149 | 3,018 | Test |
| Table | 58 | 1,108 | Test |
| Road map | 3,000 | 3,000 | Test |
| Dashboard | 73 | 1,013 | Test |
| Relation Graph | 66 | 822 | Test |
| Flowchart | 98 | 1,451 | Test |
| Visual Puzzle | 189 | 529 | Test |
| 2D Planar Layout | 25 | 252 | Test |
| **All** | 3,658 | 11,193 | Test |

Table 1: The statistics of our dataset, including eight tasks from work and life scenarios. All data were synthesized using our multi-modal self-instruct strategy.
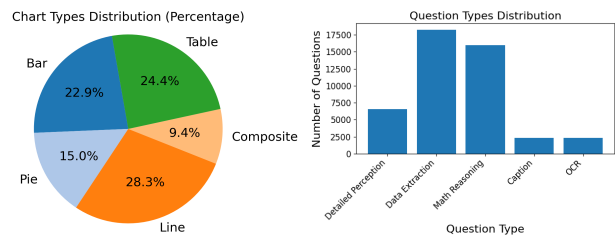


Figure 3: Left: The distribution of different chart types. Right: The number of questions for each category.

```
for an hour and a half. This means you
left at 6:40. Therefore ...
GPT-4V: 12. The clock shows the time as
1:30 ... 1:30-1.5 hours=12:00 PM ...
Gemini-1.5-pro: 1. The clock is 2:30 ... An
 hour and a half before was 1:00 ...
```

# 3 Multimodal Self-instruct Dataset

## 3.1 Dataset Statistics

We focus on eight common but under-explored scenario images, including Chart, Table, Road Map, Relation Graph, Flowchart, Visual Puzzle, Dashboard, and 2D Planar Layout. We initially synthesized a benchmark involving all 8 scenarios, containing 3,658 images and 11,193 instructions in total, to benchmark several representative LMMs. Besides, to evaluate the quality of the synthesized data, we also synthesize three training sets for chart, table, and road map tasks, comprising 34,590, 10,886, and 17,000 training instructions, respectively. As shown in Table 1, we provide detailed statistics about our synthesized dataset.

## 3.2 Synthesis Details

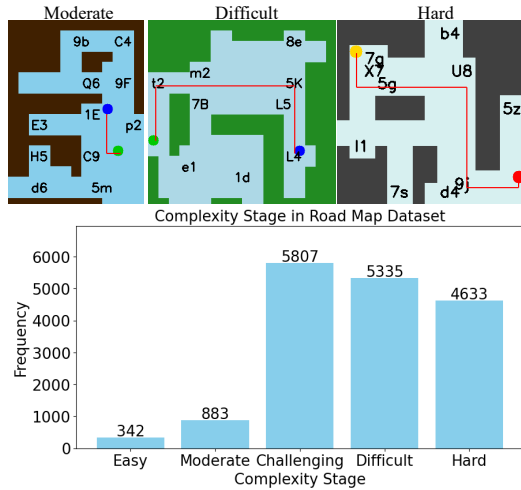**Chart and Table** Firstly, we design some keyword seeds, e.g., GDP, energy consumption, em-

Figure 4: Top: We present three examples of road maps with different path complexity. Bottom: We categorize all maps into five levels of complexity.

ployment rate, and then we prompt the LLM to expand these seed keywords into a huge keyword library covering economics, technology, and society domains. Before generation, we first randomly sample a keyword from the library and then prompt the LLM to generate corresponding visual ideas, code, and instruction data. We synthesize five types of charts: *line charts, bar charts, pie charts, table screenshots, and composite charts (containing multiple sub-charts)*. For each chart, we prompt LLMs to self-instruct five types of questions: *Optical Character Recognition (OCR), Caption, Detailed Perception (involving issues of position, quantity, layout), Data Extraction, and Mathematical Reasoning*. As shown in Figure 3, we provide statistics based on chart types and question types separately. Besides, we provide several detailed examples for each type of chart and question in Figure A2.

**Road map Navigation.** To generate simulated maps with obstacles and paths, we design a path generation strategy based on the rapidly exploring random tree algorithm[3]: Starting from an initial point, the agent randomly walks within an under-explored map, sampling the path according to the predefined walking parameters, including direction, probability, and maximum walking steps. The process stops when the maximum walking steps are reached, and the stopping position is set as the endpoint. When synthesizing maps, the LLM first sets the map size, and randomly walking parameters. Then it generates code to implement our path gen-

---

[3] https://en.wikipedia.org/wiki/Rapidly_exploring_random_tree

eration process. Ultimately, we synthesized $17k$ training maps and $3k$ testing maps. Based on the path complexity, we categorized all maps into five levels. As shown in Figure 4, most maps are of medium difficulty or higher, requiring at least two intersections and turns to reach the endpoint. We provide two complete cases in Figure A4.

**Other Scenarios Synthesis.** We employ similar processes to synthesize images of the other five scenarios, producing 1,013 Dashboard, 822 Relation Graph, 1,451 Flowchart, 529 Visual Puzzle, and 252 Planar Layout instructions. Specifically, for Flowchart, we synthesize two types: algorithm flowcharts and operating workflow. For the Relation Graph, we generate graphs with different structures, such as trees or graphs. For Dashboard, we synthesize circular dials, such as clocks, speedometers, and fuel gauges, and some elongated dials like thermometers and barometers. Regarding the Visual Puzzle task, we synthesize two types of puzzles: visual pattern induction and multi-subgraph comparison. As for the 2D Planar Layout, we synthesize architectural layouts, webpage layouts, and more. These instructions are all used as test benchmarks to evaluate the current mainstream LMMs performance. We provide some visualized cases for each task in Figures A5 to A8.

### 3.3 Implementation Details

**LLM and Prompts.** We employ *gpt-4-turbo-2024-04-09* to implement our data synthesis: idea proposal, code generation, and instruction construction. A detailed prompt is shown in Appendix A.

**Dataset Quality.** To ensure the quality of the synthesized data, we filtered the data at three levels: **code feasibility, image aesthetics, and answer accuracy**. I. If the generated code fails to run, we prompt the LLM to self-reflect based on the error feedback from the compiler. If the LLM still cannot produce valid code after three retries, we discard that visual idea. II. For each synthesized image, we employed Llava-1.5 (Liu et al., 2024a) to check the image aesthetics, including whether visual elements within the image interfere with each other, the reasonableness of the layout, and the legibility of any text. These rules allowed us to filter out aesthetically unpleasing images. III. To ensure answer accuracy, we adopted the self-consistency (Wang et al., 2022a) for answer generation: instructing the LLM to generate multiple responses based on the idea, code, and question, and then selecting the

| LMMs | Acc (%) | | |
|---|---|---|---|
| | **Chart** | **Table** | **Road Map** |
| GPT-4-Vision-1106 | **50.6** | **75.8** | 23.3 |
| Claude-3-Sonnet | 46.4 | 68.4 | 38.3 |
| Qwen-VL-Plus-70B | 40.1 | 51.6 | 18.6 |
| Vanilla Llava-1.5-7B | 10.5 | 15.8 | 0.3 |
| Vanilla Llava-1.5-13B | 13.4 | 18.3 | 5.1 |
| InstructBLIP-7B | 8.8 | 7.7 | 0.4 |
| InstructBLIP-13B | 2.8 | 2.1 | 0.6 |
| Deepseek-VL-Chat-1.3B | 18.4 | 24.2 | 9.6 |
| Deepseek-VL-Chat-7B | 25.2 | 31.1 | 18.8 |
| Llava-our-$62k$ | 30.3 $\uparrow_{19.8}$ | 51.8 $\uparrow_{36.0}$ | **67.7** $\uparrow_{67.4}$ |

Table 2: Our model is fine-tuned on chart, table, and roadmap tasks. The arrows indicate the improvements compared to Vanilla Llava-1.5-7B.

| Data Selection | Size | Chart (%) | Table (%) | Map (%) |
|---|---|---|---|---|
| Vanilla Llava | 0 | 10.5 | 15.8 | 0.3 |
| $w/$ Chart | $34.5k$ | 29.8 | 26.7 | 8.9 |
| $w/$ Table | $10.8k$ | 17.3 | 47.8 | 6.0 |
| $w/$ Map | $17k$ | 9.8 | 10.3 | 62.0 |
| $w/$ Chart, Table | $45.3k$ | 31.0 | 50.4 | 7.6 |
| $w/$ Chart, Table, Map | $62.3k$ | 30.3 | 51.8 | 67.7 |

Table 3: We investigate the synergistic effects between the three tasks. Chart and table corpus can improve each other and both benefit road map tasks.

final answer through a voting process. IV. Additionally, we randomly selected 10% of the question-answer pairs for human verification. The results confirmed that the quality of our dataset is assured.

## 4 Experiments

First, we evaluate the performance of many representative LMMs using our benchmark containing all tasks in Section 4.2. Next, we perform instruction fine-tuning on the Llava-1.5-7B using 62,476 charts, tables, and road map instructions (denoted as Llava-our-$62k$). Then, we discuss the in-domain performance Llava-our-$62k$ and the impact of the quantity of synthetic data (Section 4.3). Lastly, we investigate whether it can be generalized to other reasoning tasks (Section 4.4).

### 4.1 Settings

We evaluated the performance of mainstream open-source and closed-source LMMs, including Llava-1.5-7B (Liu et al., 2024a), Llava-1.5-13B, InstructBLIP-7B (Dai et al., 2024), InstructBLIP-13B, Deepseek-VL-Chat-1.3B (Lu et al., 2024), Deepseek-VL-Chat-7B, GPT-4-Vision-1106 (Ope-

nAI, 2023), Claude-3-Sonnet[4] and Qwen-VL-Plus (Bai et al., 2023b). All models were evaluated using the same prompts and temperature settings. We provide the evaluation metrics and other training details in Appendix A.

### 4.2 Benchmarking LMM's Visual Reasoning

As shown Figure 1, we evaluate the performance of many LMMs, Llava-our-$62k$ across eight tasks, i.e., chart, table, road map, dashboard, relation graph, flowchart, visual puzzle, and planar layout. Additionally, we invited two undergraduate students to test on our benchmark. Their scores were then averaged to represent the human-level performance. The detailed results are shown in Table A1.

**Underwhelming Abstract Image Comprehension.** We observe that for these abstract images, even advanced LMMs like GPT-4V and Claude-3 achieved only 49.5% and 50.1% accuracy on average for all tasks, leaving a significant gap to human-level performance (82.1%). Surprisingly, some tasks that seem straightforward for humans, such as planning a route on a map and recognizing clocks, prove challenging for LMMs. Specifically, in the dashboard task, even the best LMMs only achieved an accuracy of 36.2%. In the chart and relation graph tasks, we observe that LMMs often make errors when dealing with abstract concepts and spatial relationships. For example, in the Planar Layout task, GPT-4 often fails to accurately distinguish the size of the three bedrooms and whether they contain a washroom. These results indicate that despite significant progress in understanding semantic-rich natural photos, current LMMs still possess only a rudimentary understanding of abstract images and concepts.

**Significant Disparity in Visual Reasoning Ability Among LMMs.** In the road map navigation task, LMMs need to dynamically plan reasonable paths based on visual input. In the visual puzzle task, LMMs should carefully observe the given diagrams, induce visual patterns, and then perform reasoning. For these two tasks, we observed a significant performance disparity between open-source and closed-source LMMs. For example, Claude-3 achieved 38.3% and 47% for road map and visual puzzles, respectively, while smaller open-source models all achieved very low accuracy ($\leq 20\%$).

---

[4] https://www.anthropic.com/news/claude-3-family

| LLM | Weak-related Tasks (%) | | Our Synthetic Benchmark (%) | | | | |
|---|---|---|---|---|---|---|---|
| | ChartQA | MathVista | Dashboard | Relation Graph | Flowchart | Visual Puzzle | Planar Layout |
| Vanilla Llava | 19.9 | 25.1 | 16.5 | 29.6 | 9.6 | 3.4 | 37.7 |
| Llava-our-62$k$ | 23.9 ↑4 | 25.9 ↑0.8 | 16.5 | 30.1 ↑0.5 | 12.3 ↑2.7 | 3.6 ↑0.2 | 44.1 ↑6.4 |

Table 4: We used two weakly related tasks and our synthetic benchmarks from five untrained tasks to evaluate the generalization capability of our $62k$ model, which was fine-tuned solely on chart, table, and road map tasks.

This disparity between open-source and closed-source LMMs is particularly pronounced in these visual reasoning tasks.

### 4.3 Main Results After Fine-tuning

In addition to constructing the benchmark, we fine-tuned the Llava-1.5-7B model using the training sets from chart, table, and map tasks, and compared its performance with other baselines.

**In-domain Performance.** First, as shown in Table 2, compared to vanilla Llava-1.5-7B, we significantly improved its chart understanding capabilities by 19.8% and 36%, and also achieved the best performance in the road map navigation task (67.7%), far surpassing closed-source LMMs like GPT-4 (23.3%) and Claude-3 (38.3%). Notably, we only use $68k$ synthetic data and 4 hours of LoRA fine-tuning, elevating the chart understanding capability of Llava-1.5-7B to the Qwen-VL-Plus level. This demonstrates the tremendous potential of our synthetic data. Besides, we observe that most LMMs perform poorly on the road map navigation task, but can quickly improve after fine-tuning using our data. This highlights that current LMMs are not well-aligned in these reasoning scenarios.

**Synergy Between Chart, Table and Road Map.** We also studied the synergistic effects among the three tasks, such as whether chart training data benefits table and road map navigation tasks. As shown in Table 3, we trained separately on the chart ($34.5k$), table ($10.8k$), and roadmap ($17k$) datasets. Then, we train with a mix of chart and table data, and finally with a mix of all three tasks. We found that training on charts and tables does have a positive effect on road map tasks. For example, training solely on charts or tables can lead to approximately a +5% performance improvement in road map tasks, despite the significant differences in task types. Interestingly, the reverse is not true. The training process on road maps does not have a significant impact on chart and table tasks. We speculate that this may be due to the different capabilities required for each task.

**Impact of Synthetic Data Quantity.** To investigate the impact of synthetic data quantity, we fine-tuned the Llava-1.5-7B model using $35k$, $47k$, and $62k$ synthetic instructions respectively. As shown in Figure A1, we observe that as the quantity of synthetic data increases, the model's performance steadily improves without reaching a plateau, especially in the math reasoning sub-task. Specifically, the accuracy for chart tasks increased from 25.78% to 29.5%, and the table accuracy improved by 5.4%. These results indicate that our synthetic data are of high quality and diversity.

### 4.4 Generalized to Untrained Tasks

We evaluate whether Llava-our-62$k$ can generalize to other benchmarks, especially the tasks with significant differences. We use 1) two weakly correlated tasks: ChartQA (Masry et al., 2022), MathVista (Lu et al., 2023), and 2) our synthetic benchmarks from other five reasoning tasks. As shown in Table 4, we observe that although our $62k$ model is only trained on chart, table, and road map data, it also demonstrates improvements in other benchmarks, including chartQA (+4%), MathVista (+0.8%), and our synthetic benchmarks (+1.95% on average). These results show that our model can generalize to other types of visual reasoning tasks, rather than merely fitting to the training scenarios.

### 5 Conclusions

We observe that current LMMs perform sub-optimally in perceiving and reasoning with abstract images, often failing at simple daily tasks. Therefore, we design a multimodal self-instruct strategy, enabling LLMs to autonomously synthesize various diagrams, instrument dashboards, and visual puzzles using code, and self-propose reasoning Q&A. We synthesized $11k$ data to benchmark the current LMMs. Evaluation results underscore the significant challenges posed by our benchmark. We also synthesized $62k$ chart and road map training instructions to fine-tune a Llava-7B, enhancing its chart interpretation and map navigation abilities.

## Limitations

Our multi-modal strategy can synthesize a vast amount of abstract images and reasoning instructions, providing valuable training data to enhance LMMs. However, we want to highlight that there remain some limitations or areas for improvement: 1. Our data synthesis process relies on the code generation and reasoning capabilities of LLMs, which are only available in closed-source models like GPT-4. Using these models is costly. As the capabilities of open-source models improve, we are attempting to use open-source LLMs, such as Llama 3 and Deepseek-V2, to synthesize data. This will significantly reduce our expenses. 2. This work used code to synthesize abstract images in eight scenarios, such as tables and maps. In the future, we can expand to more scenarios, such as using code to control robot simulators to generate specific house layouts and structures, thereby producing a massive amount of data. 3. We believe that the image resolution of visual encoders is a bottleneck for current LMMs, especially for these abstract diagrams. In the future, we plan to improve the image resolution of the encoders to enhance the fine-grained perception capabilities of LMMs.

## References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. 2022. Flamingo: a visual language model for few-shot learning. *Preprint*, arXiv:2204.14198.

AI Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*.

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023a. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023b. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond.

Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. 2024. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*.

Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023a. Sharegpt4v: Improving large multimodal models with better captions. *arXiv preprint arXiv:2311.12793*.

Sijin Chen, Xin Chen, China. Xiaoyan Zhang, Mingsheng Li, Gang Yu, Hao Fei, Hongyuan Zhu, Jiayuan Fan, and Tao Chen. 2023b. Ll3da: Visual interactive instruction tuning for omni-3d understanding, reasoning, and planning. In *CVPR*.

Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. 2024. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in Neural Information Processing Systems*, 36.

Danny Driess, F. Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Ho Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Peter R. Florence. 2023. Palm-e: An embodied multimodal language model. In *International Conference on Machine Learning*.

Yifan Du, Hangyu Guo, Kun Zhou, Wayne Xin Zhao, Jinpeng Wang, Chuyuan Wang, Mingchen Cai, Ruihua Song, and Ji-Rong Wen. 2023. What makes for good visual instructions? synthesizing complex visual reasoning instructions for visual instruction tuning. *arXiv preprint arXiv:2311.01487*.

Yucheng Han, Chi Zhang, Xin Chen, Xu Yang, Zhibin Wang, Gang Yu, Bin Fu, and Hanwang Zhang. 2023. Chartllama: A multimodal llm for chart understanding and generation. *arXiv preprint arXiv:2311.16483*.

Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. 2023. Cogagent: A visual language model for gui agents. *arXiv preprint arXiv:2312.08914*.

Anwen Hu, Yaya Shi, Haiyang Xu, Jiabo Ye, Qinghao Ye, Ming Yan, Chenliang Li, Qi Qian, Ji Zhang, and

9

Fei Huang. 2023. mplug-paperowl: Scientific diagram analysis with the multimodal large language model. *arXiv preprint arXiv:2311.18248*.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui, Owais Khan Mohammed, Qiang Liu, Kriti Aggarwal, Zewen Chi, Johan Bjorck, Vishrav Chaudhary, Subhojit Som, Xia Song, and Furu Wei. 2023. Language is not all you need: Aligning perception with language models. *ArXiv*, abs/2302.14045.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.

Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. 2024. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. *arXiv preprint arXiv:2401.13649*.

Nicholas Lee, Thanakul Wattanawong, Sehoon Kim, Karttikeya Mangalam, Sheng Shen, Gopala Anumanchipali, Michael W Mahoney, Kurt Keutzer, and Amir Gholami. 2024. Llm2llm: Boosting llms with novel iterative data enhancement. *arXiv preprint arXiv:2403.15042*.

Bin Lei, Yuchen Li, and Qiuwu Chen. 2024. Autocoder: Enhancing code large language model with {AIEV-Instruct}.

Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. 2023a. Otter: A multi-modal model with in-context instruction tuning. *ArXiv*, abs/2305.03726.

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

Junnan Li, Dongxu Li, S. Savarese, and Steven Hoi. 2023b. Blip-2: Bootstrapping Language-Image Pretraining with Frozen Image Encoders and Large Language Models. *ArXiv*, abs/2301.12597.

Yanda Li, Chi Zhang, Gang Yu, Zhibin Wang, Bin Fu, Guosheng Lin, Chunhua Shen, Ling Chen, and Yunchao Wei. 2023c. Stablellava: Enhanced visual instruction tuning with synthesized image-dialogue data. *arXiv preprint arXiv:2308.10253*.

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024b. Visual instruction tuning. *Advances in neural information processing systems*, 36.

Ruibo Liu, Jerry Wei, Fangyu Liu, Chenglei Si, Yanzhe Zhang, Jinmeng Rao, Steven Zheng, Daiyi Peng, Diyi Yang, Denny Zhou, et al. 2024c. Best practices and lessons learned on synthetic data for language models. *arXiv preprint arXiv:2404.07503*.

Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2023. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*.

Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, et al. 2024. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*.

Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2023. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*.

Ahmed Masry, Parsa Kavehzadeh, Xuan Long Do, Enamul Hoque, and Shafiq R. Joty. 2023. Unichart: A universal vision-language pretrained model for chart comprehension and reasoning. *ArXiv*, abs/2305.14761.

Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*.

Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruti Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. 2024. Mm1: Methods, analysis & insights from multimodal llm pre-training. *arXiv preprint arXiv:2403.09611*.

OpenAI. 2022. Chatgpt.

OpenAI. 2023. Gpt-4 technical report.

Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. 2023. Kosmos-2: Grounding multimodal large language models to the world. *ArXiv*, abs/2306.14824.

Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. 2023. Pandagpt: One model to instruction-follow them all. *ArXiv*, abs/2305.16355.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aur'elien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and Efficient Foundation Language Models. *ArXiv*, abs/2302.13971.

Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models. *ArXiv*, abs/2307.09288.

Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou. 2022a. Self-consistency improves chain of thought reasoning in language models. *ArXiv*, abs/2203.11171.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022b. Self-instruct: Aligning language model with self generated instructions. *Preprint*, arXiv:2212.10560.

Ziao Wang, Yuhang Li, Junda Wu, Jaehyeon Soon, and Xiaofeng Zhang. 2023. Finvis-gpt: A multimodal large language model for financial chart analysis. *ArXiv*, abs/2308.01430.

Jingxuan Wei, Nan Xu, Guiyong Chang, Yin Luo, Bi-Hui Yu, and Ruifeng Guo. 2024. mchartqa: A universal benchmark for multimodal chart question answer based on vision-language alignment and reasoning. *arXiv preprint arXiv:2404.01548*.

Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. 2023. Magicoder: Source code is all you need. *arXiv preprint arXiv:2312.02120*.

Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2023a. Next-gpt: Any-to-any multimodal llm. *ArXiv*, abs/2309.05519.

Weijia Wu, Yuzhong Zhao, Hao Chen, Yuchao Gu, Rui Zhao, Yefei He, Hong Zhou, Mike Zheng Shou, and Chunhua Shen. 2023b. Datasetdm: Synthesizing data with perception annotations using diffusion models. *Advances in Neural Information Processing Systems*, 36:54683–54695.

Renqiu Xia, Bo Zhang, Hao Peng, Ning Liao, Peng Ye, Botian Shi, Junchi Yan, and Y. Qiao. 2023. Structchart: Perception, structuring, reasoning for visual chart understanding. *ArXiv*, abs/2309.11268.

Junlin Xie, Zhihong Chen, Ruifei Zhang, Xiang Wan, and Guanbin Li. 2024. Large multimodal agents: A survey. *arXiv preprint arXiv:2402.15116*.

Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023a. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*.

Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. 2023b. Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models. *arXiv preprint arXiv:2306.09265*.

Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Mingshi Yan, Yuhao Dan, Chenlin Zhao, Guohai Xu, Chenliang Li, Junfeng Tian, Qiang Qi, Ji Zhang, and Feiyan Huang. 2023a. mplug-docowl: Modularized multimodal large language model for document understanding. *ArXiv*, abs/2307.02499.

Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yi Zhou, Junyan Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qiang Qi, Ji Zhang, and Feiyan Huang. 2023b. mplug-owl: Modularization empowers large language models with multimodality. *ArXiv*, abs/2304.14178.

Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2023. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*.

Zhenfei Yin, Jiong Wang, Jianjian Cao, Zhelun Shi, Dingning Liu, Mukai Li, Xiaoshui Huang, Zhiyong Wang, Lu Sheng, Lei Bai, et al. 2024. Lamm: Language-assisted multi-modal instruction-tuning dataset, framework, and benchmark. *Advances in Neural Information Processing Systems*, 36.

Kaining Ying, Fanqing Meng, Jin Wang, Zhiqian Li, Han Lin, Yue Yang, Hao Zhang, Wenbo Zhang, Yuqi Lin, Shuo Liu, et al. 2024. Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi. *arXiv preprint arXiv:2404.16006*.

11

Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2023a. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*.

Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023b. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*.

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

Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. 2023. Glm-130b: An Open Bilingual Pre-trained Model. *ICLR 2023 poster*.

Jun Zhan, Junqi Dai, Jiasheng Ye, Yunhua Zhou, Dong Zhang, Zhigeng Liu, Xin Zhang, Ruibin Yuan, Ge Zhang, Linyang Li, Hang Yan, Jie Fu, Tao Gui, Tianxiang Sun, Yugang Jiang, and Xipeng Qiu. 2024. Anygpt: Unified multimodal llm with discrete sequence modeling. *ArXiv*, abs/2402.12226.

Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. 2023. Appagent: Multimodal agents as smartphone users. *Preprint*, arXiv:2312.13771.

Liang Zhang, Anwen Hu, Haiyang Xu, Mingshi Yan, Yichen Xu, Qin Jin, Ji Zhang, and Fei Huang. 2024a. Tinychart: Efficient chart understanding with visual token merging and program-of-thoughts learning. *ArXiv*, abs/2404.16635.

Wenqi Zhang, Ke Tang, Hai Wu, Mengna Wang, Yongliang Shen, Guiyang Hou, Zeqi Tan, Peng Li, Yueting Zhuang, and Weiming Lu. 2024b. Agentpro: Learning to evolve via policy-level reflection and optimization. *arXiv preprint arXiv:2402.17574*.

Henry Hengyuan Zhao, Pan Zhou, and Mike Zheng Shou. 2023. Genixer: Empowering multimodal large language models as a powerful data generator. *arXiv preprint arXiv:2312.06731*.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *ArXiv*, abs/2304.10592.

## A  Experiments Details

**Metrics.** Considering the diversity of output formats, including numerical values, single phrases, and long sentences, we employed different evaluation metrics. For numerical questions in chart, table, and dashboard tasks, answers within a 5% error margin are considered correct. For numerical questions in other tasks, the predicted values must match the labeled values exactly. For single-phrase answers, the predictions should either precisely match or contain the labeled answers. For long-sentence answers, we used the Rouge-L score as the evaluation metric. For the map navigation task, we evaluated the predicted paths by calculating the Landmark Coverage Rate (LCR(%)): we first extracted the predicted landmark sequence from the LMM's response and then compared it sequentially with the annotated landmarks sequence, calculating the proportion of correctly ordered landmarks.

**Training Details.** We fine-tuned the Llava-1.5-7B using LoRA (Hu et al., 2021) (denoted as Llava-our-62$k$) on chart, table, and road map training sets for 1 epoch, with a batch size of 16, a learning rate of 2e-4, a rank of 128 and alpha of 256. All other parameters were kept consistent with those of Llava-1.5-7B. For reasoning questions, we concatenated the answer and rationale for instruction-following training.

## B  Additional Experiment Results

As discussed in Section 4.2, we evaluate the performance of many LMMs, Llava-our-62k and humans using our benchmark. All results are shown in Table A1. Besides, as shown in Table B2, we also calculated the Rough-L score for the caption sub-task in the chart and table.

## C  Related Work

### C.1  Multi-modal LLMs

With the rapid development of Large Language Models (LLM), many researchers are currently devoting their efforts to developing multimodal large models (MLLM) for visual understanding and reasoning tasks. Beyond OpenAI's GPT-4V and Google's Gemini, numerous open-sourced MLLMs have also emerged and gained significant progress.

Recently, MLLMs commonly align visual perception with LLMs to acquire multimodal perceptions through lightweight vision-to-language adapters, including projection, Q-former and addi-
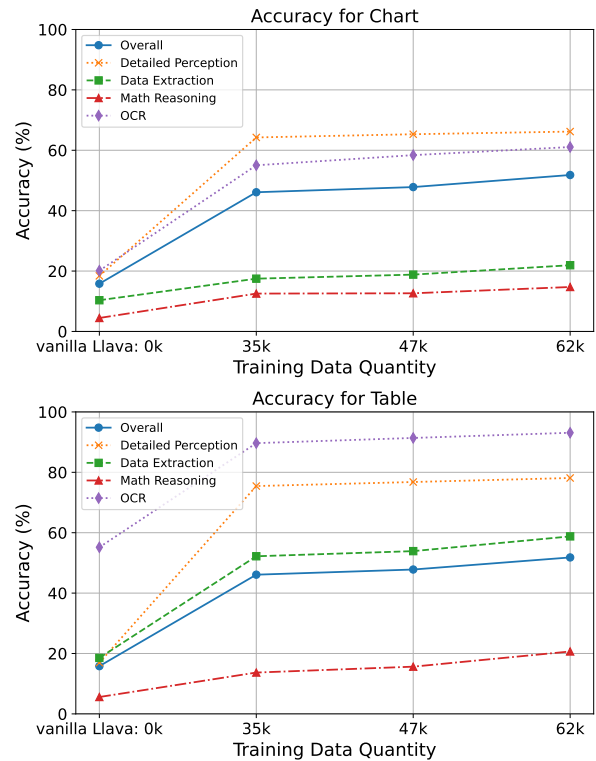


Figure A1: We analyzed the impact of synthetic data quantity on the model's performance. We fine-tune Llava-1.5-7B using chart and table instruction data of varying scales and report its accuracy. Additionally, we report the accuracy for four sub-category tasks: Detailed Perception, Data Extraction, Math Reasoning, and OCR.

tional cross-attention layers. For example, Kosmos-1/2 (Huang et al., 2023; Peng et al., 2023) and LLaVA-series models (Liu et al., 2024b,a) adopt a linear layer or an MLP to project visual inputs into textual embeddings. Furthermore, PaLM-E (Driess et al., 2023), PandaGPT (Su et al., 2023), NExT-GPT (Wu et al., 2023a) and AnyGPT (Zhan et al., 2024) even project other multimodal data such as audio, video and robot sensor data into the textual embeddings. Q-former was first proposed in BLIP-2 (Li et al., 2023b) by employing a set of learnable queries to bridge the gap between a frozen image encoder and the LLM. It has been used in several other approaches, such as LL3DA (Chen et al., 2023b), minigpt-4 (Zhu et al., 2023), Instruct-BLIP (Dai et al., 2024) and mPLUG-Owl (Ye et al., 2023b). Additionally, Flamingo (Alayrac et al., 2022) and Otter (Li et al., 2023a) inserted additional cross-attention layers into the frozen LLM to bridge the vision-only and language-only models.

However, those models are primarily focused on natural images, and there still remain challenges in

13

| LLMs | Acc (%) | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Chart | Table | Road Map | Dashboard | Relation Graph | Flowchart | Visual Puzzles | Layout | Avg. |
| Human | **93.5** | **95.1** | **75.0** | **85.3** | **82.5** | **65.5** | **62.5** | **97.6** | **82.1** |
| GPT-4Vision-1106 | 50.6[*] | 75.8[*] | 23.3 | 36.2[*] | 52.4 | 45.3[*] | 35.9 | 76.6[*] | 49.5 |
| Claude-3-Sonnet | 46.4 | 68.4 | 38.3 | 35.4 | 56.2[*] | 40.3 | 47.0[*] | 69.1 | 50.1 |
| Qwen-VL-Plus | 40.1 | 51.6 | 18.6 | 26.4 | 52.2 | 32.5 | 32.3 | 61.5 | 39.4 |
| Deepseek-VL-Chat-7B | 25.2 | 31.1 | 18.8 | 18.2 | 37.6 | 20.8 | 15.0 | 47.2 | 26.7 |
| Vanilla Llava-1.5-7B | 10.5 | 15.8 | 0.3 | 16.5 | 29.6 | 9.6 | 3.4 | 37.7 | 15.4 |
| Llava-our-62$k$ | 30.3 | 51.8 | 67.7[*] | 16.5 | 30.1 | 12.3 | 3.6 | 44.1 | 32.0 |

Table A1: Evaluating LMMs using our synthesized benchmark containing eight reasoning tasks. Bold indicates the best performance. [*] indicates the second highest.

| LLMs | Rough-L | |
| --- | --- | --- |
| | Chart | Table |
| GPT-4Vision-1106 | 0.42 | 0.42 |
| Claude-3-Sonnet | 0.48 | 0.46 |
| Qwen-VL-Plus | 0.36 | 0.37 |
| Vanilla Llava-1.5-7B | 0.33 | 0.37 |
| Vanilla Llava-1.5-13B | 0.33 | 0.40 |
| InstructBLIP-7B | 0.04 | 0.23 |
| InstructBLIP-13B | 0.05 | 0.11 |
| Deepseek-VL-Chat-1.3B | 0.36 | 0.35 |
| Deepseek-VL-Chat-7B | 0.39 | 0.37 |
| Llava-our-62k | 0.46 | 0.44 |

Table B2: For the chart and table tasks, we also calculated the captioning results.

the comprehension of complex fine-grained images such as charts, documents and diagrams. Some multimodal benchmarks have made valuable explorations into the visual reasoning capabilities and fine-grained recognition abilities of LMMs (Yue et al., 2024; Yin et al., 2024; Xu et al., 2023b; Antol et al., 2015; Liu et al., 2023; Ying et al., 2024; Yu et al., 2023b; Li et al., 2024).

Besides, several MLLMs have been proposed for chart comprehension and reasoning, including ChartLlama (Han et al., 2023), Unichart (Masry et al., 2023), Structchart (Xia et al., 2023), FinVis-GPT (Wang et al., 2023) and TinyChart (Zhang et al., 2024a). mPLUG-DocOwl (Ye et al., 2023a) strengthens the OCR-free document understanding ability with a document instruction tuning dataset.

## C.2 Data Synthesis

Data synthesis is widely used in LLM training to supplement the insufficiency of instruction-following data. Many studies focus on generating high-quality synthetic data either distilling dialogue data from a strong LLM (Wang et al., 2022b; Xu et al., 2023a; Yu et al., 2023a; Chen et al., 2023a; Zhao et al., 2023), or using external tools to refine LLM-generated synthetic data (Wei et al., 2023; Lee et al., 2024). For instance, Wang et al. (2022b) proposed *Self-Instruct* to improve the instruction-following ability of LLMs via their own generation of instruction data. Xu et al. (2023a) further generated more complex instruction through *Evol-Instruct*. Yu et al. (2023a) synthesized a mathematical dataset from LLMs by bootstrapping mathematical questions and rewriting the question from multiple perspectives. Wei et al. (2023) can generate diverse and realistic coding problems from open-source code snippets. Lei et al. (2024) can also create high-quality large code datasets for LLMs. It simulates programmers writing code and conducting unit tests through agent interactions, ensuring annotation accuracy with an external code executor.
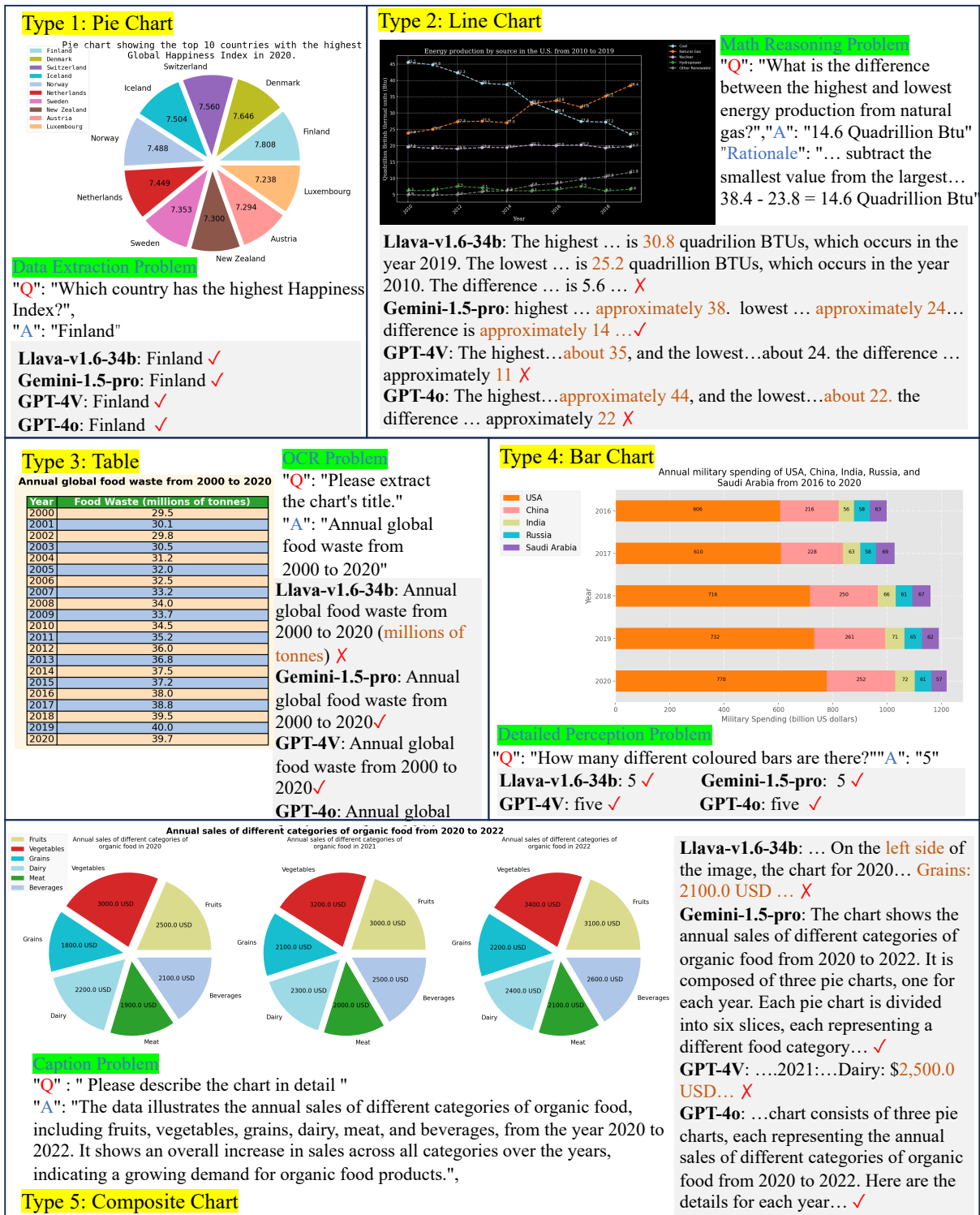
Figure A2: The chart task includes five types of charts (pie chart, line chart, table, bar chart, composite chart), each containing five types of questions (Data Extraction, Math Reasoning, OCR, Detailed Perception, Caption Problem).

**" Q " :** " You are now in a 18 * 18 road map. The road map is divided into two areas: the obstacle area and the free area. The color of obstacle area is Olive Green, and the color of free area is Light Blue. You need to find a path from the start point to the end point. The color of start point is Green, and the color end point is Blue. You can move up, down, left, or right. The road map has some markers to help you find the path. Now you need to find a path from the start point to the end point. "

**"A":** "To navigate through the road map and reach the end point, follow these steps:
1. Start at the green point.
2. Move upwards until you reach the marker labeled "t2".
3. Turn right and continue moving in that direction for a certain number of steps (not specified) until you reach the marker labeled "5K".
4. Turn downwards and continue moving in that direction for a certain number of steps (not specified) until you reach the marker labeled "L4".
5. Move right for one step and you will reach the end point, which is labeled blue."

**Llava-v1.6-34b**: Move up until … "t2" … Move right until … "m2" Move down until …"5k". Move left until…"L5"… Move down reach…"L4… ✓ Score = 1/1
**Gemini-1.5-pro**: Move right until … "m2" Move right until… "8e" Move downwards until …"5K" Move downwards until …"L5" Move downwards until "L4" … end point. ✗ Score = 0.7/1
**GTP-4V**: Move upwards until … "t2"…move right reach …"m2 " … Move downwards reach "7B" … Move right…Move downwards … reach "L4" ✗ Score = 0.7/1

**GPT-4o**: To find a path from the start point (green dot) to the end point (blue dot) in the maze, follow these steps:
1. Move right until you reach the marker labeled "e1".
2. Move up to the marker labeled "m2".
3. Move right until you reach the marker labeled "8e".
4. Move down to the marker labeled "5K".
5. Move down to the marker labeled "L5".
6. Move down to reach the end point (blue dot labeled "L4").
✗ Score = 0.7/1



**" Q ":** "You are now in a 15 by 15 road map. The road map is divided into two areas: the obstacle area and the free area. The color of obstacle area is Dark Gray, and the color of free area is Cream. You need to find a path from the start point to the end point. The color of start point is Red, and the color end point is Gold. You can move up, down, left, or right. The road map has some markers to help you find the path. Now you need to find a path from the start point to the end point."

**" A ":** "To navigate the road map, you will need to move through the cream-colored areas while avoiding the dark gray obstacles. Your starting point is marked by a red color, while your destination is marked with a gold color. Follow this path:
Move left until you reach the marker "9j", then move up until you reach marker "U8". From there, continue left until you reach the marker "X7". Then, move up until you reach marker "7g", and finally move left to reach the end point.

**Llava-v1.6-34b**: Move up until "t2"…move right until …"b4" move down until …"u8" … move left until…"5z" … move down until … "11" move right until …"7s" … move down until …"d4"…move right until …"9j"…move up until … end point
✗ Score = 0.3/1
**Gemini-1.5-pro**: Move upwards until … "5z".Move left until … "U8". Move upwards until … "b4". Move left until … "5g". Move upwards until … "x7". Move upwards until … "7g". Move left … end point.
✗ Score = 0.8/1
**GPT-4V**: Move upwards… "d4j"… move left until … "7s" … move upwards until …"I1 " … move left … "5g"…moving left…until … "b4" … move upwards … end point. ✗ Score = 0.0/1

**GPT-4o**: Let's find a path from the start point (Red) to the end point (Gold) in the 15x15 road map, following the markers as we go.
1.Start at the red point.
2.Move left to marker "5z".
3.Move up to marker "U8".
4.Move left to marker "b4".
5.Move down to marker "7g".
6.Move left to marker "X7".
7.Move down to the end point (Gold). ✗ Score = 0.5/1

Figure A3: We present two examples of road map navigation, including the synthesized simulated maps, questions, and answers.
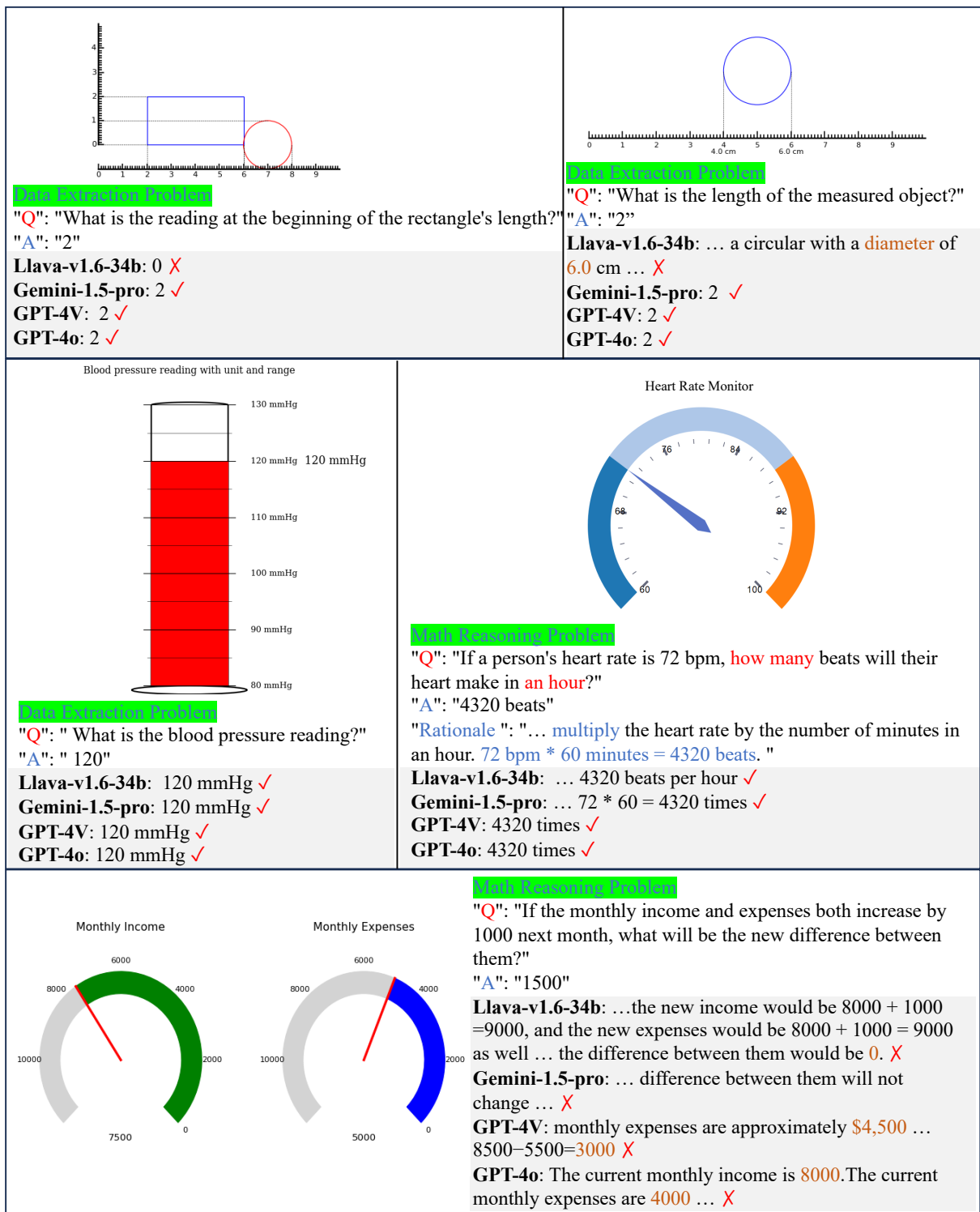
**Data Extraction Problem**
"Q": "What is the reading at the beginning of the rectangle's length?"
"A": "2"
**Llava-v1.6-34b**: 0 ✗
**Gemini-1.5-pro**: 2 ✓
**GPT-4V**: 2 ✓
**GPT-4o**: 2 ✓

**Data Extraction Problem**
"Q": "What is the length of the measured object?"
"A": "2"
**Llava-v1.6-34b**: … a circular with a diameter of 6.0 cm … ✗
**Gemini-1.5-pro**: 2 ✓
**GPT-4V**: 2 ✓
**GPT-4o**: 2 ✓

**Data Extraction Problem**
"Q": " What is the blood pressure reading?"
"A": " 120"
**Llava-v1.6-34b**: 120 mmHg ✓
**Gemini-1.5-pro**: 120 mmHg ✓
**GPT-4V**: 120 mmHg ✓
**GPT-4o**: 120 mmHg ✓

**Math Reasoning Problem**
"Q": "If a person's heart rate is 72 bpm, how many beats will their heart make in an hour?"
"A": "4320 beats"
"Rationale ": "… multiply the heart rate by the number of minutes in an hour. 72 bpm * 60 minutes = 4320 beats. "
**Llava-v1.6-34b**: … 4320 beats per hour ✓
**Gemini-1.5-pro**: … 72 * 60 = 4320 times ✓
**GPT-4V**: 4320 times ✓
**GPT-4o**: 4320 times ✓

**Math Reasoning Problem**
"Q": "If the monthly income and expenses both increase by 1000 next month, what will be the new difference between them?"
"A": "1500"
**Llava-v1.6-34b**: …the new income would be 8000 + 1000 =9000, and the new expenses would be 8000 + 1000 = 9000 as well … the difference between them would be 0. ✗
**Gemini-1.5-pro**: … difference between them will not change … ✗
**GPT-4V**: monthly expenses are approximately $4,500 … 8500−5500=3000 ✗
**GPT-4o**: The current monthly income is 8000.The current monthly expenses are 4000 … ✗

Figure A4: We present five examples of dashboard.

Figure A5: We present two examples of relation graph, each containing two types of questions.
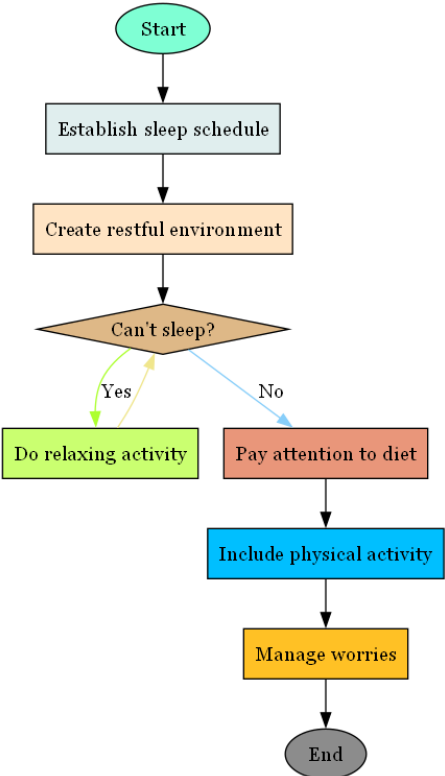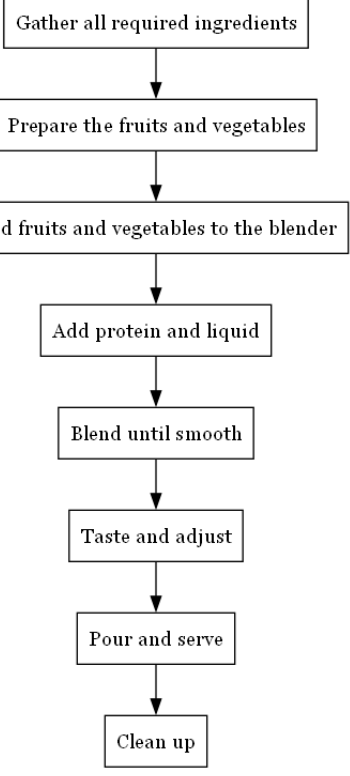
| Algorithm Workflow | Operating Workflow |
|---|---|
| Structural<br>"Q": "How many rectangular symbols are there in the flowchart?"<br>"A": "6"<br><br>**Llava-v1.6-34b**: six ✓<br>**Gemini-1.5-pro**: 5 ✗<br>**GPT-4V**: seven ✗<br>**GPT-4o**: six ✓ | Structural<br>"Q": "How many steps are there in the flowchart for preparing a healthy smoothie?"<br>"A": "8"<br><br>**Llava-v1.6-34b**: The flowchart … has six steps:1. Gather… 2. Prepare … 3. Add … 4. Add … 5. Blend … 6. Taste … 7. Pour … 8. Clean … . ✗<br>**Gemini-1.5-pro**: 8 steps ✓<br>**GPT-4V**: 8 ✓<br>**GPT-4o**: 8 ✓ |
| Reasoning<br>"Q": "What action should be taken if you can't sleep? A. Start B. Create restful environment<br>C. Do relaxing activity. D. Pay attention to diet"<br>"A": "C. Do relaxing activity"<br><br>**Llava-v1.6-34b**: Do relaxing activity ✓<br>**Gemini-1.5-pro**: Do a relaxing activity ✓<br>**GPT-4V**: Do relaxing activity ✓<br>**GPT-4o**: Do relaxing activity ✓ | Reasoning<br>"Q": "What is the previous step of 'Clean up'? A. Blend until smooth B. Taste and adjust C. Pour and serve D. None of the above"<br>"A": "C. Pour and serve"<br><br>**Llava-v1.6-34b**: Pour and serve ✓<br>**Gemini-1.5-pro**: Pour and serve ✓<br>**GPT-4V**: Pour and serve ✓<br>**GPT-4o**: Pour and serve ✓ |



Figure A6: We present two examples of flowchart (algorithm workflow and operating workflow), each containing two kinds of questions (Structural and Reasoning Problem).
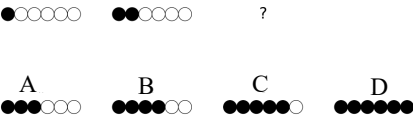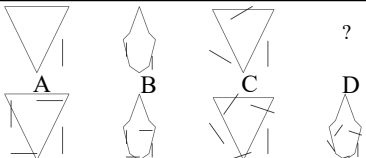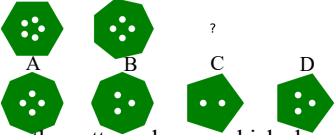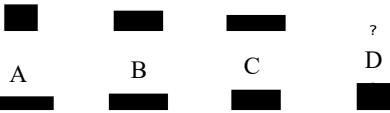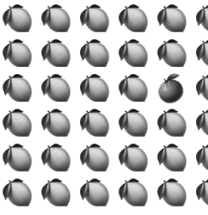
**Visual Pattern Reasoning**



"Q": "What could be the following shape in the figure ? (choose from the options) ""A": "A"

**Llava-v1.6-34b**: a square with a circle in the center … ✗
**Gemini-1.5-pro**: A ✓
**GPT-4V**: A ✓
**GPT-4o**: B, third shape should have three black circles followed by two white circles ✗

"Q": " What might be the next shape? (choose from the options) " "A": "B"

**Llava-v1.6-34b**: without more ... difficult to predict ... ✗
**Gemini-1.5-pro**: C, Column 1: Triangles with lines added in each step … ✗
**GPT-4V**: please clarify the specific pattern, so I can help…✗
**GPT-4o**: D , the lines inside the initial triangle move to form the outline of the next distorted shape ✗

"Q": Based on the pattern shown, which shape comes next? (choose from the options) ""A": "B"

**Llava-v1.6-34b**: without more clear pattern … difficult to predict … ✗
**Gemini-1.5-pro**: B ✓
**GPT-4V**: Option C, Shape B: Hexagon with 5 dots. ✗
**GPT-4o**: D… second shape is a hexagon with five dots. ✗

"Q" :"Given the pattern in the image, what is the next shape? (choose from the options ) ""A": "A"

**Llava-v1.6-34b**: appears to be a rectangle ✗
**Gemini-1.5-pro**: A ✓
**GPT-4V**: pattern shows a consistent movement of the black rectangle to the right … ✗
**GPT-4o**: D… decreasing width of the rectangle ✗

**Multi-Subgraph Comparison**

"Q": "Identify the different object by its row and column in the image. "A": "(3, 4)"

**Llava-v1.6-34b**: Row 3, Column 5: A pair of blue jeans ✗
**Gemini-1.5-pro**: 3,4,Pants ✓
**GPT-4V**: 4th row, 4th column…a pair of pants ✗
**GPT-4o**: 4th row, 3rd column ✓

"Q": " Identify the different object by its row and column. " "A": "(5, 5)"

**Llava-v1.6-34b**:challenging to identify ✗
**Gemini-1.5-pro**: same ✗
**GPT-4V**: same ✗
**GPT-4o**: (5th row, 5th column) ✓

| 3! | 11 - 5 | 9 - 4 | 1 + 5 |
| 3! | 7 - 1 | 5 + 1 | 2 + 4 |
| 6 / 1 | 6 + 0 | 3! | 4 + 2 |
| 5 + 1 | 7 - 1 | 5 + 1 | 6 / 1 |

"Q": "What different result does one equation in the image produce?" "A": "5"

**Llava-v1.6-34b**: … interpret "3!" as the factorial of 3, then it would produce a different result … ✗
**Gemini-1.5-pro**: All the equations in the image result in the number six…✗
**GPT-4V**: 6 / 1 calculates to 6, but the only division operation … ✗
**GPT-4o**: 5 ✓

"Q": " Which fruit do the identical emojis in the image depict?" "A": "lemon"

**Llava-v1.6-34b**: a pear ✗
**Gemini-1.5-pro**: lemons ✓
**GPT-4V**: peaches ✗
**GPT-4o**: peaches ✗

Figure A7: We present two categories of visual puzzles (visual pattern reasoning and muti-subgraph comparison), each containing four visual puzzle graphs, questions, and answers.
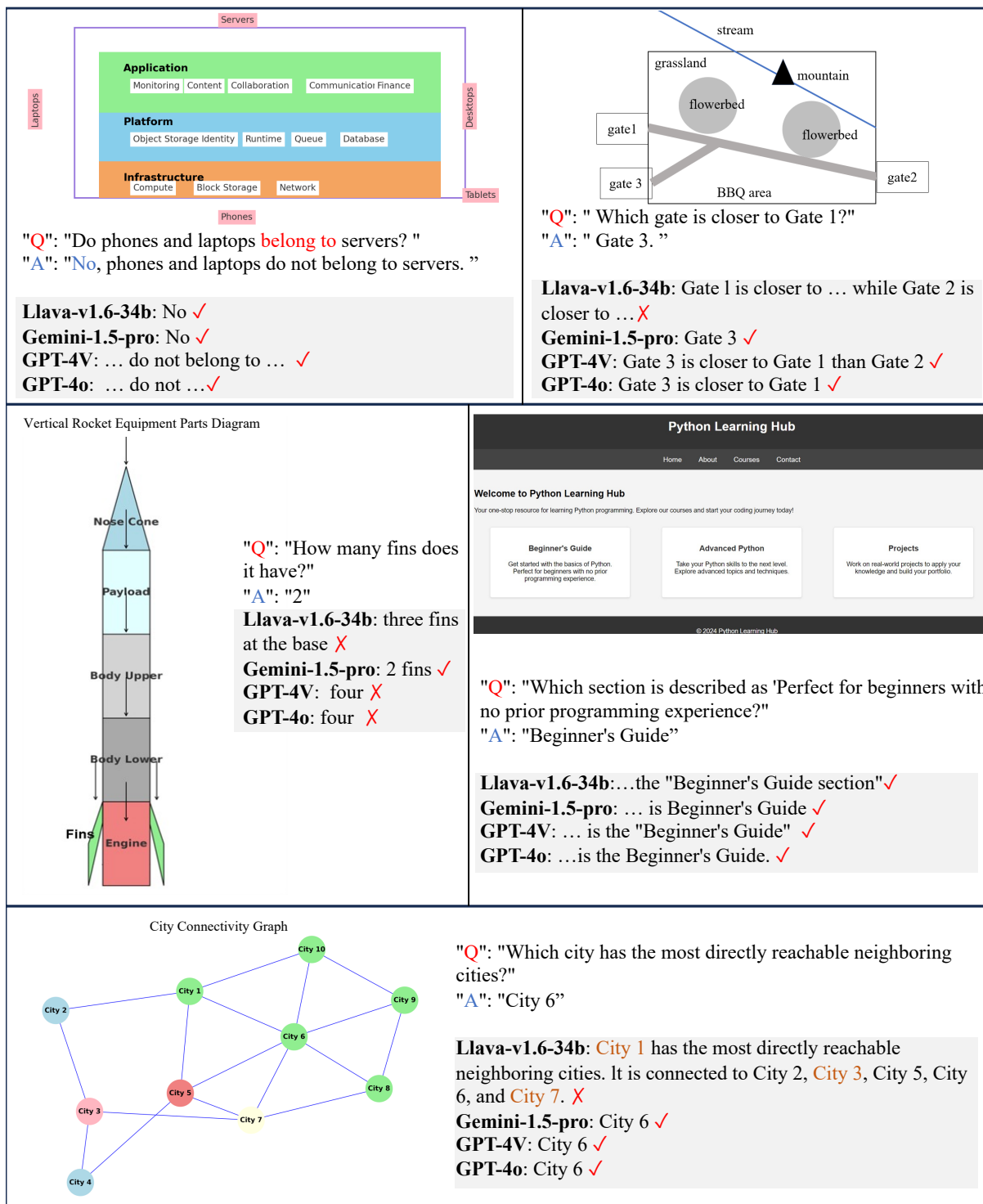
"Q": "Do phones and laptops belong to servers? "
"A": "No, phones and laptops do not belong to servers. ”

**Llava-v1.6-34b**: No ✓
**Gemini-1.5-pro**: No ✓
**GPT-4V**: … do not belong to … ✓
**GPT-4o**: … do not …✓

"Q": " Which gate is closer to Gate 1?"
"A": " Gate 3. ”

**Llava-v1.6-34b**: Gate l is closer to … while Gate 2 is closer to …✗
**Gemini-1.5-pro**: Gate 3 ✓
**GPT-4V**: Gate 3 is closer to Gate 1 than Gate 2 ✓
**GPT-4o**: Gate 3 is closer to Gate 1 ✓

Vertical Rocket Equipment Parts Diagram

"Q": "How many fins does it have?"
"A": "2"
**Llava-v1.6-34b**: three fins at the base ✗
**Gemini-1.5-pro**: 2 fins ✓
**GPT-4V**:  four ✗
**GPT-4o**: four  ✗

"Q": "Which section is described as 'Perfect for beginners with no prior programming experience?"
"A": "Beginner's Guide”

**Llava-v1.6-34b**:…the "Beginner's Guide section"✓
**Gemini-1.5-pro**: … is Beginner's Guide ✓
**GPT-4V**: … is the "Beginner's Guide" ✓
**GPT-4o**: …is the Beginner's Guide. ✓

City Connectivity Graph

"Q": "Which city has the most directly reachable neighboring cities?"
"A": "City 6”

**Llava-v1.6-34b**: City 1 has the most directly reachable neighboring cities. lt is connected to City 2, City 3, City 5, City 6, and City 7. ✗
**Gemini-1.5-pro**: City 6 ✓
**GPT-4V**: City 6 ✓
**GPT-4o**: City 6 ✓

Figure A8: We present five examples of 2D planar layout, including the layout graph, problems, answers and rationales.

```
--------------- Data Prompt ---------------
Generate data related to Digital Forensics Unit.
Requirements:
The data should describe a tree-like structure of Digital Forensics Unit.
There can be multiple layers and certain nodes can have no children.
The data should not contain too much nodes and should not be too complicated.
Increase the depth of the data, but no more than 3 nodes in the same layer.
The total number of nodes should not exceed 8.
Output format: {"data": {...}}

Instance:
{
  "data": {
    "Digital Forensics Unit": {
    "Case Management": {
      "Evidence Collection": {},
      "Analysis": {}
    },
    "Training and Development": {
      "Workshops": {},
      "Certifications": {}
    }
  }
}
}

--------------- Title Prompt ---------------
Generate a title for the data.
Requirements:
The title should be brief and concise.
The title should describe the general content of the data.
Output format: {"caption": "..." }

Instance: Digital Forensics Unit

--------------- Code Prompt ---------------
Generate high quality python code to draw a organization chart for the data.
Requirements:
The code should only use packages from ['graphviz'].
The code must conform general requirements (given in JSON format):
{
  "title": "Graphic Design Team",
  "data": [
    "all data must be used",
    "annotate the node on the organization chart"
  ],
  "layout": [
    "draw an hierarchy structured organization chart of the data",
    "nodes different levels are positioned vertically, nodes on the same level are
    positioned horizontallyuse arrows or lines to connect nodes",
    "do not show axis"
  ]
}
Output format: ```python ... ```
```
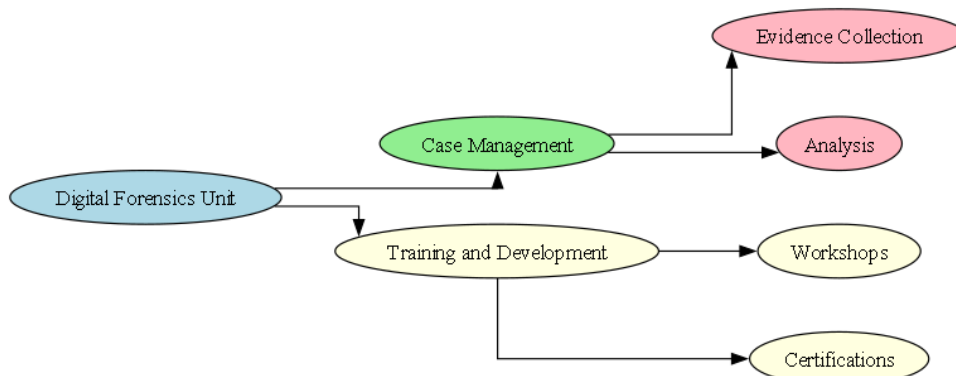
```
(continue from last page)

--------------- Question-Answer Prompt ---------------
Generate correct and high quality question-answer pairs about the data and the
organization chart.
Requirements:
Question-answer types:
{
  STRUCTURAL: {
    'Example 1': 'What is the type of this figure? Choose your answer from
    organization chart, pie chart, line chart, gantt chart.',
    'Example 2': "What's the color of {node}?"},
  MATH_REASONING: {
    'Example 1': 'Does {name} node exist in this figure?',
    'Example 2': 'How many nodes are there?'}
}
If applicable, the answer can be a single word.
Consider the data and code together to get the answer.
Output format: {
    "STRUCTURAL":[{"Q":"...", "A":"..."}, ...],
    "MATH_REASONING":[{"Q":"...", "A":"..."}, ...]
}

Instance:
{
    "STRUCTURAL": [
      {
        "Q": "What is the type of this figure? Choose your answer from
        organization chart, pie chart, line chart, gantt chart.",
        "A": "organization chart"
      },
      {
        "Q": "What's the color of the 'Digital Forensics Unit' node?",
        "A": "lightblue"
      }
    ],
    "MATH_REASONING": [
      {
        "Q": "How many nodes are there in the 'Digital Forensics Unit'?",
        "A": "2"
      },
      {
        "Q": "Does the 'Evidence Collection' node exist in this figure?",
        "A": "Yes"
      },
      {
        "Q": "How many nodes are there in the 'Case Management' department?",
        "A": "2"
      },
      {
        "Q": "How many nodes are there in the 'Training and Development'
        department?",
        "A": "2"
      },
      {
        "Q": "How many departments are there in the 'Digital Forensics Unit'?",
        "A": "2"
      }
    }
}
```