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

Wenqi Zhang^{1,*}, Zhenglin Cheng^{1,*}, Yuanyu He¹, Mengna Wang², Yongliang Shen¹
Zeqi Tan¹, Guiyang Hou¹, Mingqian He¹, Yanna Ma³, Weiming Lu^{1,†}, Yueting Zhuang¹

¹College of Computer Science and Technology, Zhejiang University

²Institute of Software, Chinese Academy of Sciences

³University of Shanghai for Science and Technology

{zhangwenqi, luwm}@zju.edu.cn

Project Page: <https://multi-modal-self-instruct.github.io>

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 pipeline, 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 Claude-3.5-Sonnet and GPT-4o 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 is available at: <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 encode all modalities empowers LMMs to understand human daily environments and execute complex tasks [Hong et al., 2023, Zhang et al., 2023b, Hu et al., 2023, Zhang et al., 2023a, 2024c, Koh et al., 2024, Zhang et al., 2024d]. 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

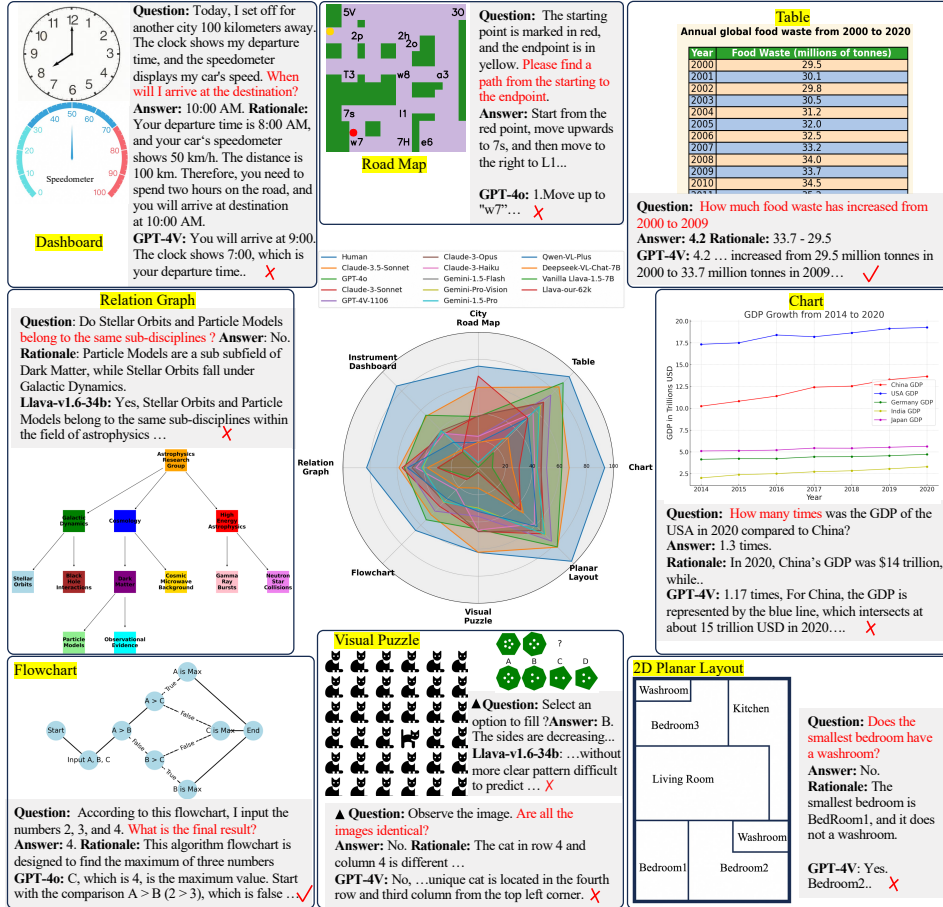


Figure 1: We leverage LLM and code to synthesize abstract images and self-instruct diverse reasoning instructions, e.g., charts, road maps, dashboards, visual puzzles, and relation graphs. Unlike natural landscapes and human photos, these non-natural images constructed with geometric elements require 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.

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., 2023b, 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, 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 [Betker et al., 2023, Esser et al., 2024], 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. After 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 rationales.

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-4o) only achieved a score of 54.7, far below the human level of 85.3. Our abstract image benchmark further indicates that the gap between current open-source models and closed-source models remains significant, despite their comparable performance on semantics-related benchmarks.

Besides, to verify the quality of the synthesized data, we 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 multi-modal 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,

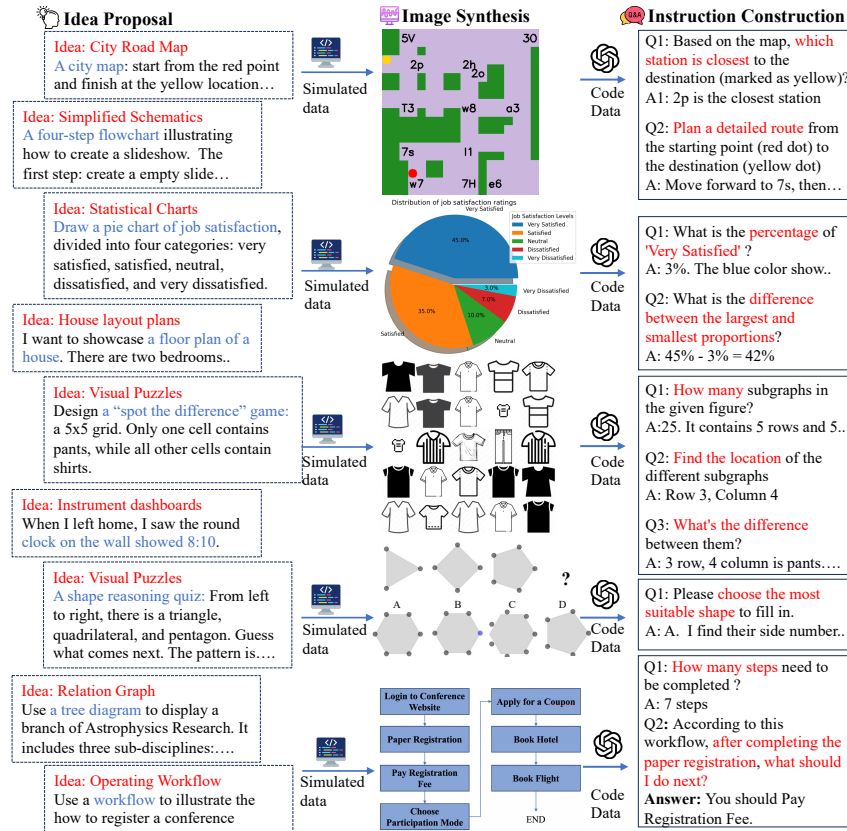


Figure 2: Our multi-modal self-instruct strategy first self-proposes a visual idea to depict an abstract image. Based on this, the LLM generates simulated data and writes code to create the drawings. Subsequently, LLM is instructed to design multiple Q&A based on the code and idea, covering various aspects such as spatial reasoning, color recognition, and mathematical reasoning, constructing a rich set of multimodal instructions.

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: 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 information. 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., Matplotlib¹ or ECharts², 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.

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
```

¹<https://matplotlib.org>

²<https://echarts.apache.org/zh/index.html>

Table 1: Left: The statistics of our dataset, including eight tasks from work and life scenarios. All data were synthesized using our multi-modal self-instruct strategy. Right: Detail Statistics on different request types and workflow structures. Right Our model is fine-tuned on chart, table, and roadmap tasks. The arrows indicate the improvements compared to Vanilla Llava-1.5-7B.

Task	#Image	# Instruction	#Usage	LMMs			
				Chart	Table	Map	
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				
				Acc (%)			
				Chart	Table	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

Rationale: I read the time from the clock as 8:10, and you have been exercising for an hour and a half. This means you left at 6:40...
GPT-4V: 12. The clock shows the time as 1:30.. 1:30-1.5 hours=12:00...
Gemini-1.5-pro: 1. The clock is 2:30 ... An hour and a half before..

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, employment 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 A1, 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 A4.

Road map Navigation. To generate simulated maps with obstacles and paths, we design a path generation strategy based on the rapidly exploring random tree algorithm³: 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

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

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 generation 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 A2, most maps are of medium difficulty or higher, requiring at least two intersections and turns to reach the endpoint. We provide two cases in Figure A6.

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 2D 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 LLMs performance. We provide some cases for each task in Figures A7 to A10.

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 Diversity. Firstly, in the data synthesis process, we control the generated topic of the image with many pre-defined keywords. For example, before synthesizing the chart, we designed a keyword library (e.g., GDP, energy, and employment rate) that includes various keywords from different domains covering economics, technology, and society. This strategy can control the generated content and avoid deviations. Similarly, during the image and question synthesis process, we use few-shot examples and templates to control the types of questions and images generated. For example, we generate five types of charts (bar, table, line, pie, composite) for the chart task, and also 5 types of questions (perception, extraction, math, caption, OCR). We also generate the difficulty levels of synthesized maps. The quantity for each category can be predefined in advance.

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 displeasing 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 final answer through a voting process.

Human Evaluation we also conduct a manual evaluation of the dataset. First, we randomly sampled 10% of the <question, answer> pairs from our benchmarks and invited 4 graduate students in the computer science field for manual evaluation. For each sample, we designed four evaluation criteria: **Image Aesthetics, Question Rationality, Answer Accuracy, and Image-Instruction Relevance**. The criteria for Image Aesthetics and Answer Accuracy are scored from 1 to 5 (5 being the highest), while Question Rationality and Image-Instruction Relevance are divided into three levels 1, 3, 5. The scoring criteria for each dimension and the final results of the human evaluation are shown in Appendix C and Table C2.

4 Experiments

First, we evaluate the performance of many leading LLMs 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

Table 2: We investigate the synergistic effects between the three tasks. Chart and table corpus can improve each other and both benefit road map tasks.

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 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.

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-62k	23.9 \uparrow ₄	25.9 \uparrow _{0.8}	16.5	30.1 \uparrow _{0.5}	12.3 \uparrow _{2.7}	3.6 \uparrow _{0.2}	44.1 \uparrow _{6.4}

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, Claude-3.5-Sonnet, Claude-3-Sonnet⁴, GPT-4o, GPT-4-Vision-1106 [OpenAI, 2023], Gemini-1.5-pro⁵ 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-4o and Claude-3.5-Sonnet achieved only 64.7% and 59.9% 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 54.79% (GPT-4o). 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-4v often fails to distinguish the size of the three bedrooms accurately 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

⁴<https://www.anthropic.com/news/claude-3-family>

⁵<https://deepmind.google/technologies/gemini/pro/>

reasoning. For these two tasks, we observed a significant performance disparity between open-source and closed-source LLMs. For example, Claude-3.5-Sonnet achieved 59.2% and 62.3% for road map and visual puzzles, respectively, while smaller open-source models all achieved very low accuracy ($\leq 20\%$). This disparity between open-source and closed-source LLMs 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 1, 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 LLMs 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 LLMs perform poorly on the road map navigation task, but can quickly improve after fine-tuning using our data. This highlights that current LLMs 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 2, 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 A3, 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-62k 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 3, 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 LLMs 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 LLMs. 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.

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. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- 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. Flamingo: a visual language model for few-shot learning, 2022.
- AI Anthropic. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*, 2024.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023a.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. 2023b.
- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.
- Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiu Shi Du, Zhe Fu, et al. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*, 2024.
- Lei Chen, Feng Yan, Yujie Zhong, Shaoxiang Chen, Zequn Jie, and Lin Ma. Mindbench: A comprehensive benchmark for mind map structure recognition and analysis. *arXiv preprint arXiv:2407.02842*, 2024.
- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023a.
- Sijin Chen, Xin Chen, China. Xiaoyan Zhang, Mingsheng Li, Gang Yu, Hao Fei, Hongyuan Zhu, Jiayuan Fan, and Tao Chen. Ll3da: Visual interactive instruction tuning for omni-3d understanding, reasoning, and planning. In *CVPR*, 2023b.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Danny Driess, F. Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzan 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. Palm-e: An embodied multimodal language model. In *International Conference on Machine Learning*, 2023. URL <https://api.semanticscholar.org/CorpusID:257364842>.
- Yifan Du, Hangyu Guo, Kun Zhou, Wayne Xin Zhao, Jinpeng Wang, Chuyuan Wang, Mingchen Cai, Ruihua Song, and Ji-Rong Wen. What makes for good visual instructions? synthesizing complex visual reasoning instructions for visual instruction tuning. *arXiv preprint arXiv:2311.01487*, 2023.

- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, and Rongrong Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models. *ArXiv*, abs/2306.13394, 2023. URL <https://api.semanticscholar.org/CorpusID:259243928>.
- Yucheng Han, Chi Zhang, Xin Chen, Xu Yang, Zhibin Wang, Gang Yu, Bin Fu, and Hanwang Zhang. Chartllama: A multimodal llm for chart understanding and generation. *arXiv preprint arXiv:2311.16483*, 2023.
- Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. Cogagent: A visual language model for gui agents. *arXiv preprint arXiv:2312.08914*, 2023.
- Anwen Hu, Yaya Shi, Haiyang Xu, Jiabo Ye, Qinghao Ye, Ming Yan, Chenliang Li, Qi Qian, Ji Zhang, and Fei Huang. mplug-paperowl: Scientific diagram analysis with the multimodal large language model. *arXiv preprint arXiv:2311.18248*, 2023.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- 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. Language is not all you need: Aligning perception with language models. *ArXiv*, abs/2302.14045, 2023. URL <https://api.semanticscholar.org/CorpusID:257219775>.
- Zhen Huang, Zengzhi Wang, Shijie Xia, Xuefeng Li, Haoyang Zou, Ruijie Xu, Run-Ze Fan, Lyuman-shan Ye, Ethan Chern, Yixin Ye, et al. Olympicarena: Benchmarking multi-discipline cognitive reasoning for superintelligent ai. *arXiv preprint arXiv:2406.12753*, 2024.
- 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. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- Yoonsik Kim, Moonbin Yim, and Ka Yeon Song. Tablevqa-bench: A visual question answering benchmark on multiple table domains. *arXiv preprint arXiv:2404.19205*, 2024.
- Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. *arXiv preprint arXiv:2401.13649*, 2024.
- Nicholas Lee, Thanakul Wattanawong, Sehoon Kim, Karttikeya Mangalam, Sheng Shen, Gopala Anumanchipali, Michael W Mahoney, Kurt Keutzer, and Amir Gholami. Llm2llm: Boosting llms with novel iterative data enhancement. *arXiv preprint arXiv:2403.15042*, 2024.
- Bin Lei, Yuchen Li, and Qiuwu Chen. Autocoder: Enhancing code large language model with {AIEV-Instruct}. 2024. URL <https://api.semanticscholar.org/CorpusID:270045303>.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *ArXiv*, abs/2305.03726, 2023a. URL <https://api.semanticscholar.org/CorpusID:258547300>.
- Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. Seed-bench: Benchmarking multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13299–13308, 2024.

- Junnan Li, Dongxu Li, S. Savarese, and Steven Hoi. Blip-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models. *ArXiv*, abs/2301.12597, 2023b.
- Yanda Li, Chi Zhang, Gang Yu, Zhibin Wang, Bin Fu, Guosheng Lin, Chunhua Shen, Ling Chen, and Yunchao Wei. Stablellava: Enhanced visual instruction tuning with synthesized image-dialogue data. *arXiv preprint arXiv:2308.10253*, 2023c.
- Fuxiao Liu, Xiaoyang Wang, Wenlin Yao, Jianshu Chen, Kaiqiang Song, Sangwoo Cho, Yaser Yacoob, and Dong Yu. Mmc: Advancing multimodal chart understanding with large-scale instruction tuning. *arXiv preprint arXiv:2311.10774*, 2023a.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2024a.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024b.
- Ruibo Liu, Jerry Wei, Fangyu Liu, Chenglei Si, Yanzhe Zhang, Jinmeng Rao, Steven Zheng, Daiyi Peng, Diyi Yang, Denny Zhou, et al. Best practices and lessons learned on synthetic data for language models. *arXiv preprint arXiv:2404.07503*, 2024c.
- Yexin Liu, Zhengyang Liang, Yueze Wang, Muyang He, Jian Li, and Bo Zhao. Seeing clearly, answering incorrectly: A multimodal robustness benchmark for evaluating mllms on leading questions. *arXiv preprint arXiv:2406.10638*, 2024d.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*, 2023b.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022.
- Ahmed Masry, Parsa Kavehzadeh, Xuan Long Do, Enamul Hoque, and Shafiq R. Joty. Unichart: A universal vision-language pretrained model for chart comprehension and reasoning. *ArXiv*, abs/2305.14761, 2023. URL <https://api.semanticscholar.org/CorpusID:258865561>.
- Ahmed Masry, Mehrad Shahmohammadi, Md Rizwan Parvez, Enamul Hoque, and Shafiq Joty. Chartinstruct: Instruction tuning for chart comprehension and reasoning. *arXiv preprint arXiv:2403.09028*, 2024.
- Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruvi Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. Mm1: Methods, analysis & insights from multimodal llm pre-training. *arXiv preprint arXiv:2403.09611*, 2024.
- Fanqing Meng, Wenqi Shao, Quanfeng Lu, Peng Gao, Kaipeng Zhang, Yu Qiao, and Ping Luo. Chartassisstant: A universal chart multimodal language model via chart-to-table pre-training and multitask instruction tuning. *arXiv preprint arXiv:2401.02384*, 2024.
- OpenAI. Chatgpt. 2022.
- OpenAI. Gpt-4 technical report. 2023.
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *ArXiv*, abs/2306.14824, 2023. URL <https://api.semanticscholar.org/CorpusID:259262263>.

- Shubhankar Singh, Purvi Chaurasia, Yerram Varun, Pranshu Pandya, Vatsal Gupta, Vivek Gupta, and Dan Roth. Flowvqa: Mapping multimodal logic in visual question answering with flowcharts. *arXiv preprint arXiv:2406.19237*, 2024.
- Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. Pandagpt: One model to instruction-follow them all. *ArXiv*, abs/2305.16355, 2023. URL <https://api.semanticscholar.org/CorpusID:258947721>.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9568–9578, 2024.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and Efficient Foundation Language Models. *ArXiv*, abs/2302.13971, 2023a.
- 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. Llama 2: Open foundation and fine-tuned chat models. *ArXiv*, abs/2307.09288, 2023b. URL <https://api.semanticscholar.org/CorpusID:259950998>.
- Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, and Neel Joshi. Is a picture worth a thousand words? delving into spatial reasoning for vision language models. *arXiv preprint arXiv:2406.14852*, 2024a.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Huai hsin Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *ArXiv*, abs/2203.11171, 2022a.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions, 2022b.
- Ziao Wang, Yuhang Li, Junda Wu, Jaehyeon Soon, and Xiaofeng Zhang. Finvis-gpt: A multimodal large language model for financial chart analysis. *ArXiv*, abs/2308.01430, 2023. URL <https://api.semanticscholar.org/CorpusID:260438486>.
- Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, et al. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. *arXiv preprint arXiv:2406.18521*, 2024b.
- Jingxuan Wei, Nan Xu, Guiyong Chang, Yin Luo, BiHui Yu, and Ruifeng Guo. mchartqa: A universal benchmark for multimodal chart question answer based on vision-language alignment and reasoning. *arXiv preprint arXiv:2404.01548*, 2024.
- Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Source code is all you need. *arXiv preprint arXiv:2312.02120*, 2023.

- Qiucheng Wu, Handong Zhao, Michael Saxon, Trung Bui, William Yang Wang, Yang Zhang, and Shiyu Chang. Vsp: Assessing the dual challenges of perception and reasoning in spatial planning tasks for vlms. *arXiv preprint arXiv:2407.01863*, 2024.
- Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal llm. *ArXiv*, abs/2309.05519, 2023a. URL <https://api.semanticscholar.org/CorpusID:261696650>.
- Weijia Wu, Yuzhong Zhao, Hao Chen, Yuchao Gu, Rui Zhao, Yefei He, Hong Zhou, Mike Zheng Shou, and Chunhua Shen. Datasetdm: Synthesizing data with perception annotations using diffusion models. *Advances in Neural Information Processing Systems*, 36:54683–54695, 2023b.
- Renqiu Xia, Bo Zhang, Hao Peng, Ning Liao, Peng Ye, Botian Shi, Junchi Yan, and Y. Qiao. Structchart: Perception, structuring, reasoning for visual chart understanding. *ArXiv*, abs/2309.11268, 2023. URL <https://api.semanticscholar.org/CorpusID:262067829>.
- Renqiu Xia, Bo Zhang, Hancheng Ye, Xiangchao Yan, Qi Liu, Hongbin Zhou, Zijun Chen, Min Dou, Botian Shi, Junchi Yan, et al. Chartx & chartvlm: A versatile benchmark and foundation model for complicated chart reasoning. *arXiv preprint arXiv:2402.12185*, 2024.
- Junlin Xie, Zhihong Chen, Ruifei Zhang, Xiang Wan, and Guanbin Li. Large multimodal agents: A survey. *arXiv preprint arXiv:2402.15116*, 2024.
- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*, 2023a.
- Peng Xu, Wenqi Shao, Kaipeng Zhang, Peng Gao, Shuo Liu, Meng Lei, Fanqing Meng, Siyuan Huang, Yu Qiao, and Ping Luo. Lvlm-ehub: A comprehensive evaluation benchmark for large vision-language models. *arXiv preprint arXiv:2306.09265*, 2023b.
- 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. mplug-docowl: Modularized multimodal large language model for document understanding. *ArXiv*, abs/2307.02499, 2023a. URL <https://api.semanticscholar.org/CorpusID:259360848>.
- 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. mplug-owl: Modularization empowers large language models with multimodality. *ArXiv*, abs/2304.14178, 2023b. URL <https://api.semanticscholar.org/CorpusID:258352455>.
- Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549*, 2023.
- Zhenfei Yin, Jiong Wang, Jianjian Cao, Zhelun Shi, Dingning Liu, Mukai Li, Xiaoshui Huang, Zhiyong Wang, Lu Sheng, Lei Bai, et al. Lamm: Language-assisted multi-modal instruction-tuning dataset, framework, and benchmark. *Advances in Neural Information Processing Systems*, 36, 2024.
- Kaining Ying, Fanqing Meng, Jin Wang, Zhiqian Li, Han Lin, Yue Yang, Hao Zhang, Wenbo Zhang, Yuqi Lin, Shuo Liu, et al. Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi. *arXiv preprint arXiv:2404.16006*, 2024.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023a.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023b.

- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024.
- 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. Glm-130b: An Open Bilingual Pre-trained Model. *ICLR 2023 poster*, 2023.
- 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. Anygpt: Unified multimodal llm with discrete sequence modeling. *ArXiv*, abs/2402.12226, 2024. URL <https://api.semanticscholar.org/CorpusID:267750101>.
- Chi Zhang, Zhao Yang, Jiakuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users, 2023a.
- Liang Zhang, Anwen Hu, Haiyang Xu, Mingshi Yan, Yichen Xu, Qin Jin, Ji Zhang, and Fei Huang. Tinychart: Efficient chart understanding with visual token merging and program-of-thoughts learning. *ArXiv*, abs/2404.16635, 2024a. URL <https://api.semanticscholar.org/CorpusID:269362640>.
- Tianyu Zhang, Suyuchen Wang, Lu Li, Ge Zhang, Perouz Taslakian, Sai Rajeswar, Jie Fu, Bang Liu, and Yoshua Bengio. Vcr: Visual caption restoration. *arXiv preprint arXiv:2406.06462*, 2024b.
- Wenqi Zhang, Yongliang Shen, Weiming Lu, and Yueting Zhuang. Data-copilot: Bridging billions of data and humans with autonomous workflow. *arXiv preprint arXiv:2306.07209*, 2023b.
- Wenqi Zhang, Yongliang Shen, Linjuan Wu, Qiying Peng, Jun Wang, Yueting Zhuang, and Weiming Lu. Self-contrast: Better reflection through inconsistent solving perspectives. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3602–3622, Bangkok, Thailand, 2024c. Association for Computational Linguistics. URL <https://aclanthology.org/2024.acl-long.197>.
- Wenqi Zhang, Ke Tang, Hai Wu, Mengna Wang, Yongliang Shen, Guiyang Hou, Zeqi Tan, Peng Li, Yueting Zhuang, and Weiming Lu. Agent-pro: Learning to evolve via policy-level reflection and optimization. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5348–5375, Bangkok, Thailand, 2024d. Association for Computational Linguistics. URL <https://aclanthology.org/2024.acl-long.292>.
- Henry Hengyuan Zhao, Pan Zhou, and Mike Zheng Shou. Genixer: Empowering multimodal large language models as a powerful data generator. *arXiv preprint arXiv:2312.06731*, 2023.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *ArXiv*, abs/2304.10592, 2023. URL <https://api.semanticscholar.org/CorpusID:258291930>.

References

References follow the acknowledgments in the camera-ready paper. Use unnumbered first-level heading for the references. Any choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font size to small (9 point) when listing the references. Note that the Reference section does not count towards the page limit.

[1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp. 609–616. Cambridge, MA: MIT Press.

[2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the GEneral NEural Simulation System*. New York: TELOS/Springer-Verlag.

[3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.

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-62k) 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.

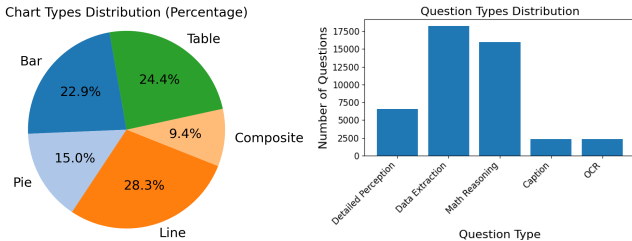


Figure A1: Left: The distribution of five chart types. Right: The number of questions for each type.

LLMs	Acc (%)								Avg.
	Chart	Table	Road Map	Dashboard	Graph	Flowchart	Puzzles	Layout	
Human	93.5	95.1	75.0	85.3	82.5	65.5	62.5	97.6	82.1
Claude-3.5-Sonnet	67.24*	84.38	59.24	54.00	58.52*	49.21	62.38*	82.94*	64.74*
GPT-4o	61.83	88.76*	37.82	54.79*	54.50	54.31*	45.37	82.54	59.99
Claude-3-Sonnet	46.4	68.4	38.3	35.4	56.2	40.3	47.0	69.1	50.1
GPT-4V-1106	50.6	75.8	23.3	36.2	52.4	45.3	35.9	76.6	49.5
GPT-4o-mini	48.7	77.4	26.7	46.3	51.1	42.5	30.8	75.8	49.5
Claude-3-Opus	46.73	67.71	38.26	38.70	48.78	35.77	47.26	65.48	48.59
internvl-2-8b	50.3	73.9	27.9	28.9	61.3	41.2	23.4	66.6	46.7
Claude-3-Haiku	41.83	57.33	23.17	35.83	45.99	23.09	45.94	58.73	41.49
Gemini-1.5-Flash	43.61	64.06	3.71	39.04	42.09	36.03	30.81	69.72	41.13
glm-4v-9b	47.8	70.9	4.4	34.3	47.0	39.3	20.2	63.8	41.0
Gemini-Pro-Vision	43.11	64.92	3.76	38.87	41.12	36.09	29.68	70.12	40.96
Gemini-1.5-Pro	43.41	63.78	3.77	38.71	41.85	35.55	30.62	69.32	40.88
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-62k	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.

B Discussion

B.1 How to improve the abstract image comprehension capabilities of LMMs?

More than just for instruction fine-tuning, our experiments have revealed that the abstract image comprehension capabilities of Large LMMs can be enhanced through various approaches:

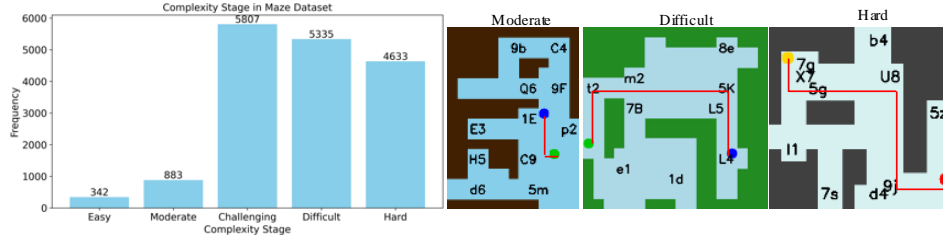


Figure A2: Left: We categorize all maps into five levels of complexity. Right: We present three examples of road maps with different path complexity.

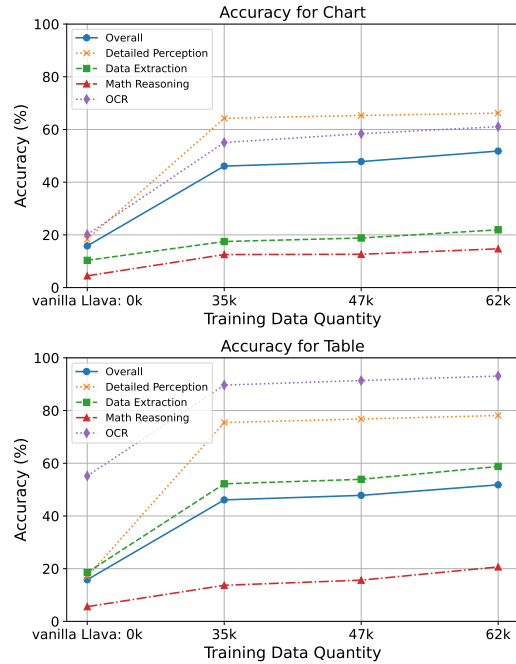


Figure A3: We analyzed the impact of synthetic data quantity on the model’s performance. We fine-tune Llava-1.5-7B using chart and table instruction data of varying scales and report its accuracy. Additionally, we report the accuracy for four sub-category tasks: Detailed Perception, Data Extraction, Math Reasoning, and OCR.

Designing More Versatile Visual Encoders: First, we observe that current LMMs have weak visual representation abilities for abstract images, which may be caused by the current visual encoders. Most of them use clip-based encoders, which emphasize semantic features while neglecting purely visual features. We plan to explore adding another visual encoder, such as DinoV2, to Llava to improve its understanding of abstract images.

Increasing Image Resolution: Then we also observed that most current LMMs resize the original image to a resolution of 336x336 when training, as it reduces the number of visual tokens without losing much semantic information. However, for these abstract images composed of lines and geometric shapes, lowering the resolution results in the loss of a significant amount of geometric features, thereby affecting the ability of abstract visual perception. Increasing the image resolution and training LMMs from scratch may be a good solution.

Incorporate into Pre-training: Next, we will incorporate abstract image data during the pre-training stage, not just during the SFT stage. This is because, while benchmarking current LMMs, we found that their weakness lies in the perception of abstract images rather than in their instruction-following abilities.

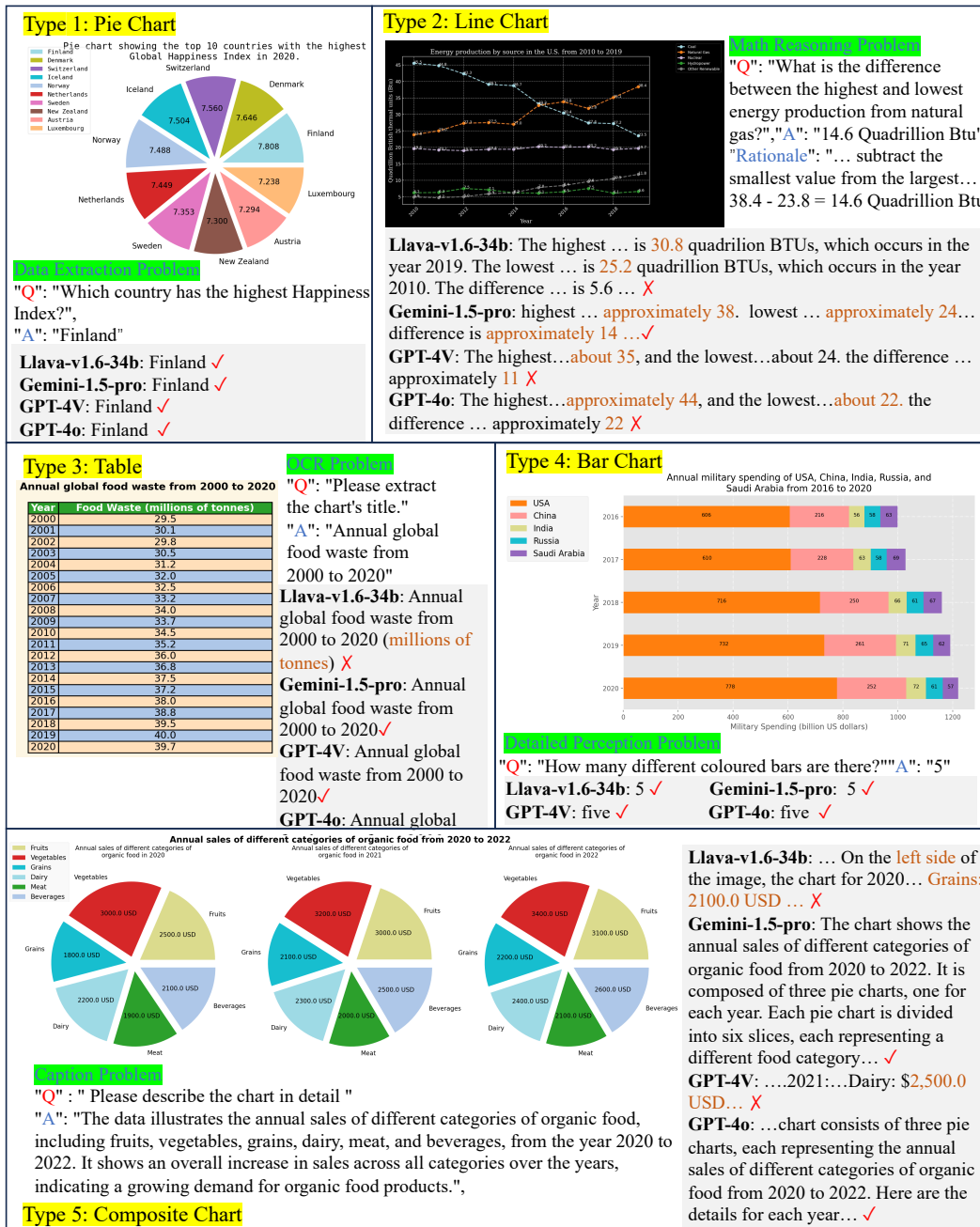


Figure A4: The chart task includes five types of charts (pie chart, line chart, table, bar chart, composite chart), each containing five types of questions (Data Extraction, Math Reasoning, OCR, Detailed Perception, Caption Problem).

Investigate Relationships between Abstract Image Tasks: Lastly, we will investigate the relationships between different abstract image tasks, quantitatively analyzing their mutual influences and their impact on LMM's abilities such as abstract image perception, spatial reasoning, and visual-symbol Induction. These fine-grained studies will guide us in designing more useful abstract image tasks using our pipeline.

We believe that focusing on the aforementioned four aspects can effectively enhance the abstract image comprehension capabilities of LMMs without compromising their general capabilities. We treat this as a direction for future work.

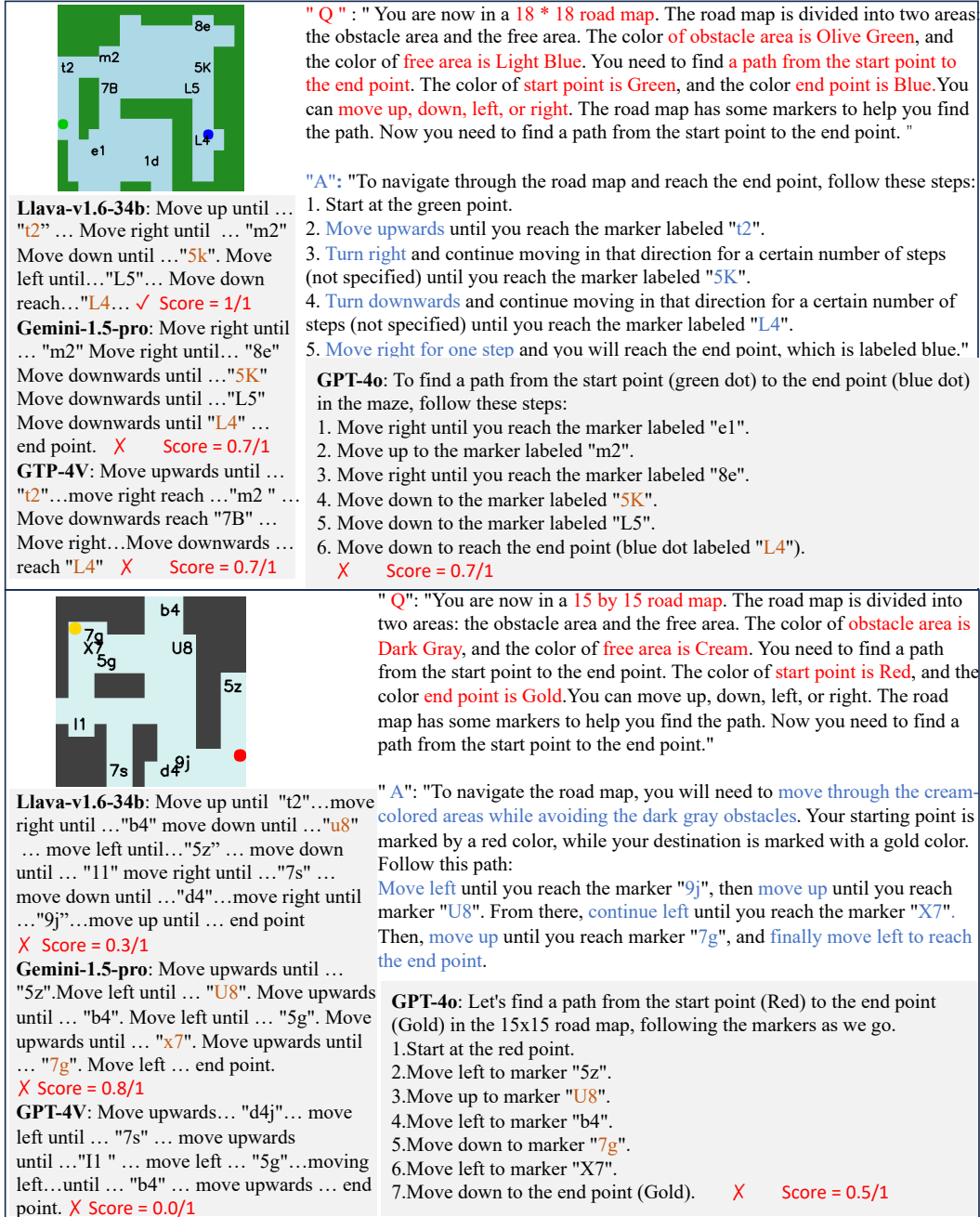


Figure A5: We present two examples of road map navigation, including the synthesized simulated maps, questions, and answers.

C Human Evaluation

As discussed in the paper, we design four evaluation metrics to manually assess the quality of the benchmark: Image Aesthetics, Question Rationality, Answer Accuracy, and Image-Instruction Relevance. The specific criteria are as follows:

- Image Aesthetics: Are the colors appropriate, are the details clearly visible, is the spatial layout reasonable, and are there any obstructions between objects?

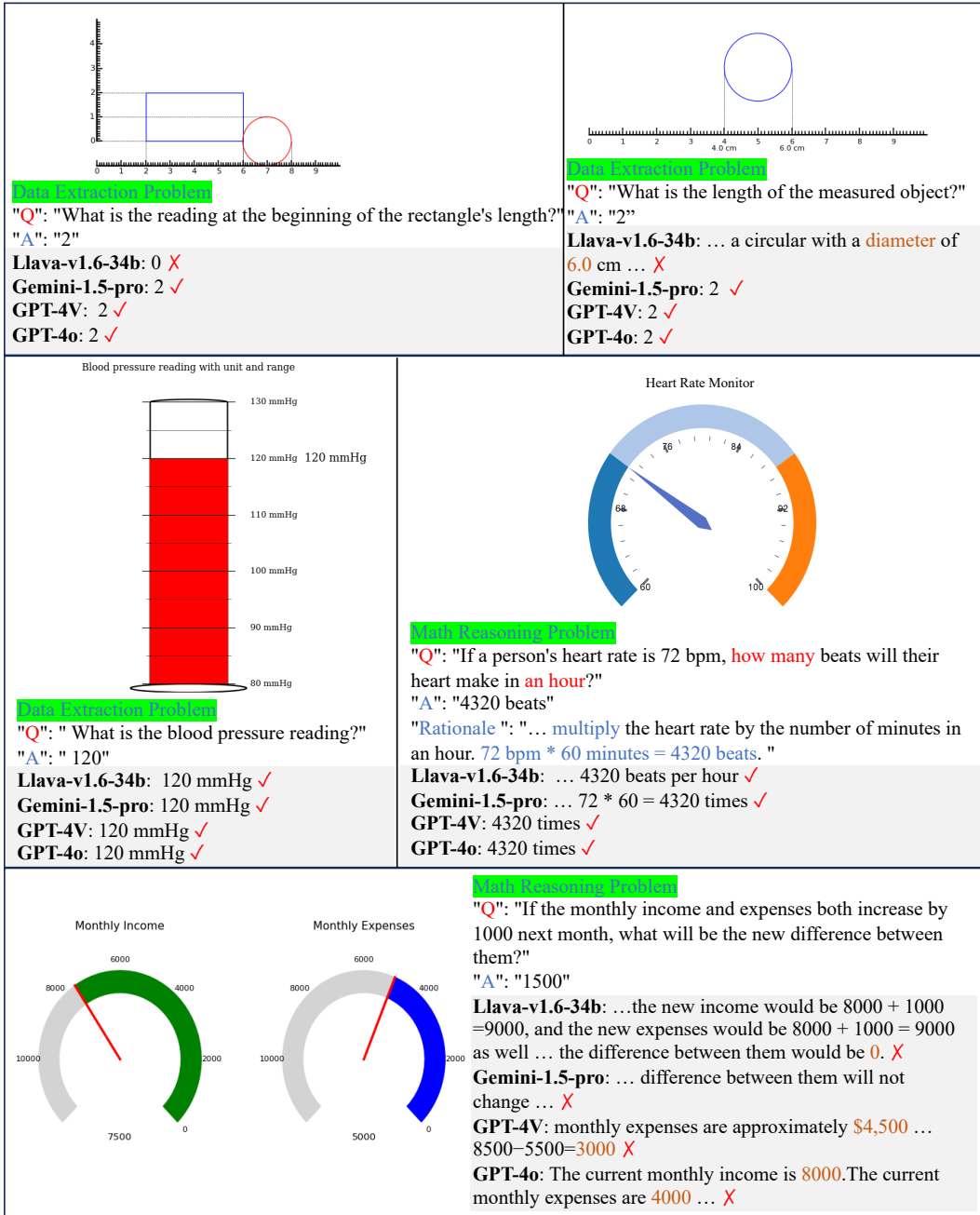


Figure A6: We present five examples of dashboard.

- Question Rationality: Is the question reasonable? Is the wording clear and unambiguous? Is the difficulty level moderate, neither too simple nor too difficult?
- Answer Accuracy: Is the answer accurate? Is the rationale logical? Are the answer and rationale consistent with each other?
- Image-Instruction Relevance: Is the answer related to the image? Can the question be answered without looking at the image?

We evaluated the benchmark both before and after filtering, with the results presented in Table C2. These results indicate that the quality of our benchmark has significantly improved post-filtering, particularly in terms of Image Aesthetics and Answer Accuracy.

Tree-based	
<p>Structural Problem</p> <p>"Q": "What is the type of this figure? Choose your answer from organization chart, pie chart, line chart, gantt chart." "A": "organization chart"</p> <p>Llava-v1.6-34b: organization chart ✓ Gemini-1.5-pro: organization chart ✓ GPT-4V: organization chart ✓ GPT-4o: organization chart ✓</p>	<p>Task Reasoning Problem</p> <p>"Q": "How many departments are there in the 'Corporate Communications Office'?", "A": "3"</p> <p>Llava-v1.6-34b: ... four departments: 1. Media Relations. 2. Internal ... 3. Crisis Management 4. Emergency Response Team ✗ Gemini-1.5-pro: ... 3 departments ... ✓ GPT-4V: three departments ... ✓ GPT-4o: three main departments ... ✓</p>
<p>Structural Problem</p> <p>"Q": "Does the 'ProjectC' node exist in this figure?", "A": "no"</p> <p>Llava-v1.6-34b: Yes, the 'ProjectC' node exists in this figure. It is connected to the 'Subproject1' node... ✗ Gemini-1.5-pro: No ✓ GPT-4V: not include ✓ GPT-4o: No ✓</p>	<p>Task Reasoning Problem</p> <p>"Q": "How many nodes are there under the 'ProjectA' node?" "A": "2"</p> <p>Llava-v1.6-34b: two nodes ✓ Gemini-1.5-pro: two nodes ✓ GPT-4V: two nodes ✓ GPT-4o: two nodes ✓</p>

Figure A7: We present two examples of relation graph, each containing two types of questions.

	Image Aesthetics	Question Rationality	Answer Accuracy	Image-Instruction Relevance
Before Filtering	2.4	3.9	3.5	4.5
After Filtering	4.0	4.1	4.3	4.4

Table C2: The results of the human evaluation.

D 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 D3, we also calculated the Rough-L score for the caption sub-task in the chart and table.

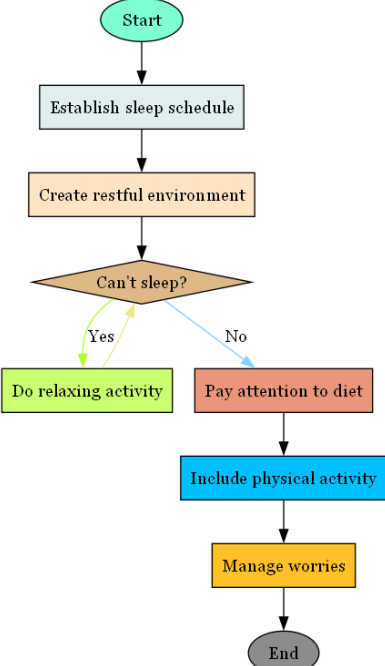
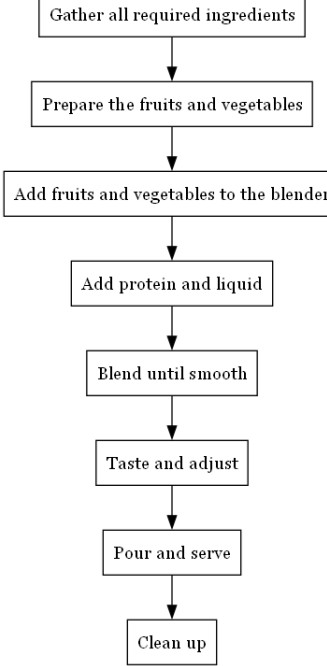
Algorithm Workflow	Operating Workflow
<p>Structural Problem</p> <p>"Q": "How many rectangular symbols are there in the flowchart?" "A": "6"</p> <p>Llava-v1.6-34b: six ✓ Gemini-1.5-pro: 5 ✗ GPT-4V: seven ✗ GPT-4o: six ✓</p>	<p>Structural Problem</p> <p>"Q": "How many steps are there in the flowchart for preparing a healthy smoothie?" "A": "8"</p> <p>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 ... ✗ Gemini-1.5-pro: 8 steps ✓ GPT-4V: 8 ✓ GPT-4o: 8 ✓</p>
<p>Reasoning Problem</p> <p>"Q": "What action should be taken if you can't sleep? A. Start B. Create restful environment C. Do relaxing activity. D. Pay attention to diet" "A": "C. Do relaxing activity"</p> <p>Llava-v1.6-34b: Do relaxing activity ✓ Gemini-1.5-pro: Do a relaxing activity ✓ GPT-4V: Do relaxing activity ✓ GPT-4o: Do relaxing activity ✓</p>	<p>Reasoning Problem</p> <p>"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" "A": "C. Pour and serve"</p> <p>Llava-v1.6-34b: Pour and serve ✓ Gemini-1.5-pro: Pour and serve ✓ GPT-4V: Pour and serve ✓ GPT-4o: Pour and serve ✓</p>
 <pre> graph TD Start([Start]) --> A[Establish sleep schedule] A --> B[Create restful environment] B --> C{Can't sleep?} C -- Yes --> D[Do relaxing activity] C -- No --> E[Pay attention to diet] E --> F[Include physical activity] F --> G[Manage worries] G --> End([End]) </pre>	 <pre> graph TD A[Gather all required ingredients] --> B[Prepare the fruits and vegetables] B --> C[Add fruits and vegetables to the blender] C --> D[Add protein and liquid] D --> E[Blend until smooth] E --> F[Taste and adjust] F --> G[Pour and serve] G --> H[Clean up] </pre>

Figure A8: We present two examples of flowchart (algorithm workflow and operating workflow), each containing two kinds of questions (Structural and Reasoning Problem).

E Related Work

E.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.

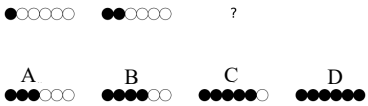
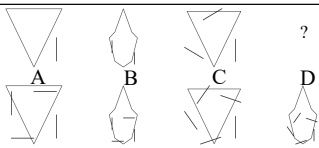
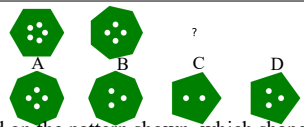


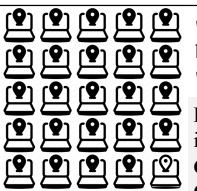

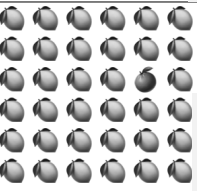
Visual Pattern Reasoning	
 <p>"Q": "What could be the following shape in the figure ? (choose from the options) ""A": "A" Llava-v1.6-34b: a square with a circle in the center ... X Gemini-1.5-pro: A ✓ GPT-4V: A ✓ GPT-4o: B, third shape should have three black circles followed by two white circles X</p>	 <p>"Q": " What might be the next shape? (choose from the options) " "A": "B" Llava-v1.6-34b: without more ... difficult to predict ... X Gemini-1.5-pro: C, Column 1: Triangles with lines added in each step ... X GPT-4V: please clarify the specific pattern, so I can help... X GPT-4o: D , the lines inside the initial triangle move to form the outline of the next distorted shape X</p>
 <p>"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 ... X Gemini-1.5-pro: B ✓ GPT-4V: Option C, Shape B: Hexagon with 5 dots. X GPT-4o: D... second shape is a hexagon with five dots. X</p>	 <p>"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 X Gemini-1.5-pro: A ✓ GPT-4V: pattern shows a consistent movement of the black rectangle to the right ... X GPT-4o: D... decreasing width of the rectangle X</p>
Multi-Subgraph Comparison	
 <p>"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 X Gemini-1.5-pro: 3,4,Pants ✓ GPT-4V: 4th row, 4th column...a pair of pants X GPT-4o: 4th row, 3rd column ✓</p>	 <p>"Q": " Identify the different object by its row and column. " "A": "(5, 5)" Llava-v1.6-34b:challenging to identify X Gemini-1.5-pro: same X GPT-4V: same X GPT-4o: (5th row, 5th column) ✓</p>
 <p>"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 ... X Gemini-1.5-pro: All the equations in the image result in the number six... X GPT-4V: 6 / 1 calculates to 6, but the only division operation ... X GPT-4o: 5 ✓</p>	 <p>"Q": " Which fruit do the identical emojis in the image depict?" "A": "lemon" Llava-v1.6-34b: a pear X Gemini-1.5-pro: lemons ✓ GPT-4V: peaches X GPT-4o: peaches X</p>

Figure A9: We present two categories of visual puzzles (visual pattern reasoning and multi-subgraph comparison), each containing four visual puzzle graphs, questions, and answers.

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

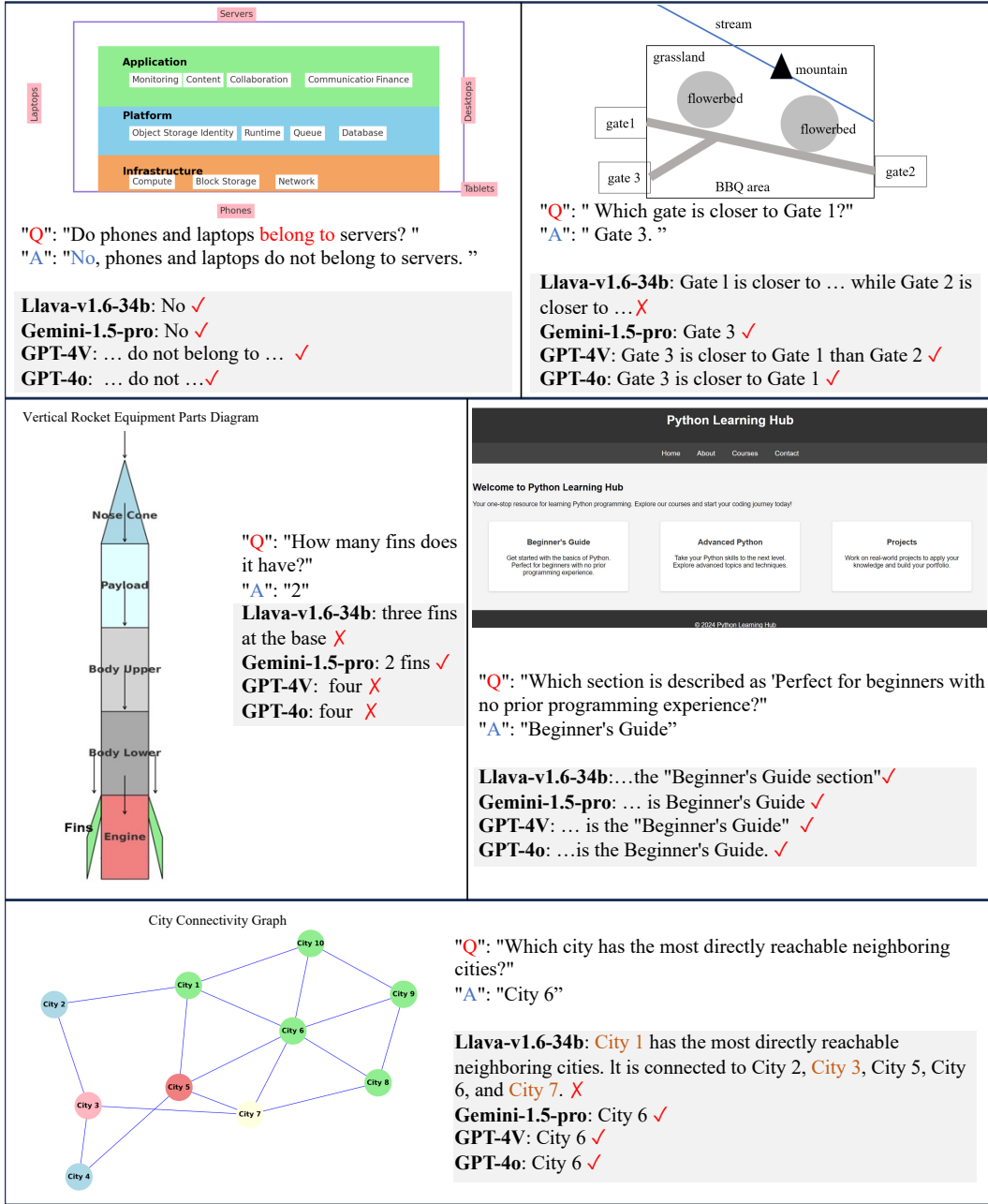


Figure A10: We present five examples of 2D planar layout, including the layout graph, problems, answers and rationales.

2023b], minigt-4 [Zhu et al., 2023], InstructBLIP [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 the comprehension of complex fine-grained images such as charts, documents, and diagrams.

----- 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": "..."}

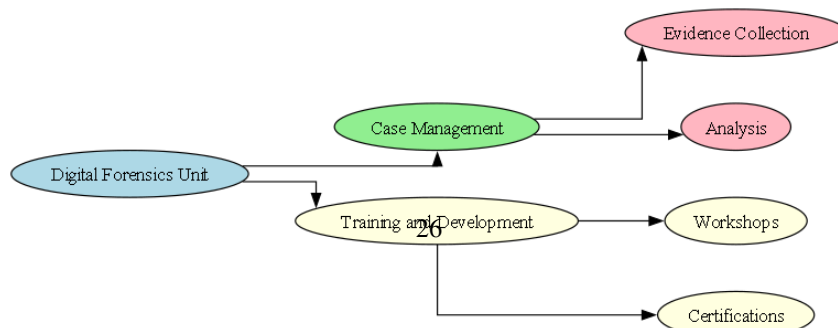
Instance: Digital Forensics Unit

----- 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 horizontally",
    "use 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"
    }
  ]
}
```

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 D3: For the chart and table tasks, we also calculated the captioning results.

E.2 Benchmark For Multimodal Model

Designing a fair benchmark to evaluate the capabilities of multimodal models has garnered widespread attention within the academic community [Antol et al., 2015, Fu et al., 2023, Xu et al., 2023b, Liu et al., 2023a, Yu et al., 2023b, Yue et al., 2024, Liu et al., 2024d, Tong et al., 2024, Huang et al., 2024]. Recently, some multimodal benchmarks have made valuable explorations into the visual reasoning capabilities and fine-grained recognition abilities of LMMs [Yin et al., 2024, Liu et al., 2023b, Ying et al., 2024, Li et al., 2024, Wang et al., 2024a, Chen et al., 2024, Wu et al., 2024, Singh et al., 2024, Zhang et al., 2024b].

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], TinyChart [Zhang et al., 2024a], CharXiv [Wang et al., 2024b], ChartX [Xia et al., 2024], TableVQA-Bench [Kim et al., 2024] and mChartQA [Wei et al., 2024]. mPLUG-DocOwl [Ye et al., 2023a] strengthens the OCR-free document understanding ability with a document instruction tuning dataset. Chartassisstant [Meng et al., 2024] undergoes a two-stage training process, starting with pre-training on chart-to-table parsing to align chart and text, followed by multitask instruction-following fine-tuning. ChartInstruct [Masry et al., 2024] employs a two-step approach to extract chart data tables and input them into the LLM. These efforts have all contributed to the advancement of multimodal technologies.

E.3 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.