

Cantor[✂]: Inspiring Multimodal Chain-of-Thought of MLLM

Timin Gao^{1*}, Peixian Chen^{2*}, Mengdan Zhang^{2*}, Chaoyou Fu², Yunhang Shen², Yan Zhang^{1†},
Shengchuan Zhang¹, Xiwu Zheng¹, Xing Sun², Liujuan Cao¹, Rongrong Ji¹

¹Key Laboratory of Multimedia Trusted Perception and Efficient Computing,

Ministry of Education of China, Xiamen University ²Tencent Youtu Lab

timinggao@stu.xmu.edu.cn, {pxchen13, zhangmengdanrz, bzhy986}@gmail.com

Abstract

With the advent of large language models (LLMs) enhanced by the chain-of-thought (CoT) methodology, visual reasoning problem is usually decomposed into manageable sub-tasks and tackled sequentially with various external tools. However, such a paradigm faces the challenge of the potential “determining hallucinations” in decision-making due to insufficient visual information and the limitation of low-level perception tools that fail to provide abstract summaries necessary for comprehensive reasoning. We argue that converging visual context acquisition and logical reasoning is pivotal for tackling visual reasoning tasks. This paper delves into the realm of multimodal CoT to solve intricate visual reasoning tasks with multimodal large language models (MLLMs) and their cognitive capability. To this end, we propose an innovative multimodal CoT framework, termed Cantor, characterized by a perception-decision architecture. Cantor first acts as a decision generator and integrates visual inputs to analyze the image and problem, ensuring a closer alignment with the actual context. Furthermore, Cantor leverages the advanced cognitive functions of MLLMs to perform as multifaceted experts for deriving higher-level information, enhancing the CoT generation process. Our extensive experiments demonstrate the efficacy of the proposed framework, showing significant improvements in multimodal CoT performance across two complex visual reasoning datasets, without necessitating fine-tuning or ground-truth rationales. Project Page: <https://ggg0919.github.io/cantor/>.

1. Introduction

With the development of large language models (LLMs), researchers have begun to adopt the chain-of-thought (CoT) strategy to improve the model performance in reason tasks.

*Equal contribution.

†Corresponding author.

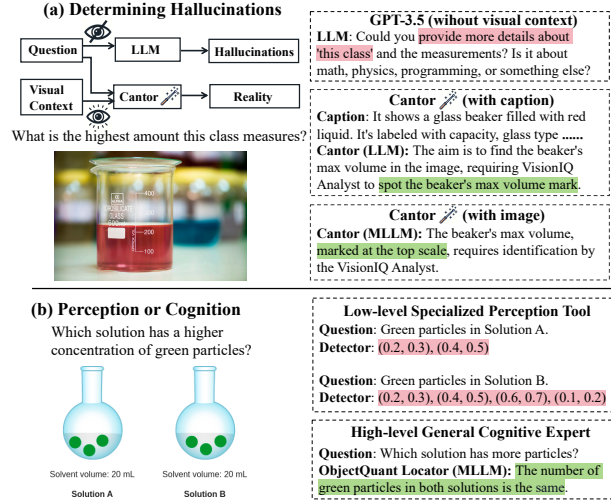


Figure 1. (a) Comparison of visual information on decision generation: Asking GPT-3.5 (without visual context) leads to “determining hallucinations” due to lacking clarity of the image. Cantor (with caption) by introducing visual context through captions, does not encounter this issue. Cantor (with image) is even more precise, improving the rationality of task assignment. (b) Comparison of different visual tools: Low-level specialized perception tools used in traditional approaches only obtain basic data. High-level general cognitive expert acted by MLLM obtains object number relationships, enabling direct and subsequent reasoning.

CoT mimics the gradual reasoning process of humans, helping models improve their deep understanding and analytical abilities by constructing a series of logical steps to solve complex visual reasoning problems. The effectiveness of CoT has been widely validated in language reasoning tasks. Recently, researchers have naturally extended its application to multimodal domains. Visual reasoning tasks [29, 30] are inherently suited for chain-of-thought (CoT) methodologies. These tasks necessitate that models not only “perceive” the contents and contexts within images but also “comprehend” these visual elements to make coherent in-

ferences and decisions. Consequently, the exploration of multimodal CoT has significantly expanded in the research community.

Most existing multimodal CoT methods are divided into two stages: decision-generation and execution. 1) **Decision-Generation**. It is the first step in multimodal CoT methods, which involves understanding, analyzing, and formulating inference plans for the problem. The existing determining methods include breaking down problems into sub-problems [53], capturing scene maps in images [32], finding similarities and differences in related images [49], and so on [41, 44]. They attempt to simplify the problem at the textual level or add more contextual information at the visual level. 2) **Execution**. In this stage, models perform specific operations scheduled by the previous determining stage. Specifically, the model transforms the planning into practical solutions. The existing execution methods usually rely on various specialized API tools or vision-language models (VLMs), with the former emphasizing the specificity of task execution [31, 41] and the latter emphasizing the universality of task execution [44, 53].

Although these multimodal CoT methods have improved the performance in visual reasoning tasks, there are still limitations: Firstly, when making decisions, existing methods often directly input plain text into LLMs without considering visual context [17, 44, 53]. Intuitively, this increases the divergent thinking of LLMs towards problems, but in reality, it may lead to “determining hallucinations”. As shown in Fig. 1 (a), if the question itself is not closely related to the image and only asks “What is the highest amount this class measures?” based on the text, LLM (GPT-3.5) is not clear about what “this class” specifically means. It will answer that the provided information is insufficient and begin to guess whether the “class” refers to a metric in physics or a class in programming. This perception uncertainty may lead LLMs to make decisions that are unrelated to the problem or even incorrect, misleading subsequent execution and resulting in completely unrelated answers.

Secondly, during execution, existing methods typically execute tasks by calling external tools, because MLLMs still fall short of solving numerous visual reasoning tasks [17, 31, 32, 38, 44]. But these tools are mostly low-level visual-perception tools (detectors, recognizer, OCR, etc.) that can only extract low-level visual information. As shown in Fig. 1 (b), when comparing the number of particles in solutions, they only provide the positions of particles and fail to infer high-level information such as the relationship between their numbers. They further input these low-level clues into LLMs for organization and summarization [17, 32, 53]. When complex clues increase, this undoubtedly increases the burden of LLMs on long-text reasoning. Meanwhile, with many external tools, it also increases the complexity of the pipeline.

To address the above limitations, we propose a novel multimodal CoT framework, Cantor. In decision generation, we enable an MLLM or an LLM to act as a cantor within the chorus, simultaneously processing visual and textual context for comprehensive understanding, and then assigning specific tasks to “experts” acted by a single MLLM for high-level logical problem-solving. Specifically, during the decision generation, we analyze in detail the importance of visual information in the determining stage. This includes the quality of determining with or without visual information, as well as the differences in the impact of detailed or concise visual information on determining. Ultimately, we conclude that visual information is crucial during the decision generation stage. When we use an MLLM model (such as Gemini) for the decision generator, we directly feed images into the model to fully comprehend the question and deliberate on it. However, when employing an LLM model (such as GPT-3.5), we find that providing a more detailed caption of the image is more conducive to understanding the question. Furthermore, the decision generator is required to explicitly provide explanatory decisions, including problem-solving strategies, reasons for expert invocation, and specific task conduction for each expert. Consequently, it guides an MLLM to act as tailored experts (such as ObjectQuant Locator, TextIntel Extractor, VisionIQ Analyst, and ChartSense Expert) to provide conclusive answers for sub-tasks in the process. As shown in Fig. 1 (a), when using LLM to make a decision, with detailed caption guidance, the model knows that it is asking for the maximum volume of the beaker and makes the correct decision. The decision is clearer when the image is available to the MLLM, that is, requiring the VisionIQ Analyst to extract the number at the top of the cup wall. During execution, we observe that MLLM is an advanced cognitive tool that performs better in directly acquiring high-level information (e.g., relative position and quantity) than acquiring low-level visual information like detecting positions. Such high-level information is superior for multimodal CoT. Instead of using several external tools, Cantor assigns different tasks to a single MLLM via different expert identities and task instructions, exploring the professional potential of an MLLM acting as certain experts. The tailored experts provide high-level professional information directly, thus reducing the burden of subsequent integrated reasoning. As shown in Fig. 1 (b), when comparing the concentration of green particles, we need to compare the number of particles in the two bottles first. MLLM acts as an ObjectQuant Locator and directly compares the quantity variance in the two solutions. Compared with obtaining the position of particles, MLLM gets the result of the quantity relationship more accurately. This result is directly applied to the further inference of the final answer.

Our proposed framework Cantor achieves SOTA results

in both ScinceQA [29] and Mathvista [30]. When Gemini is used as the decision generator, Cantor obtains an accuracy gain of 4.11% and 5.9%, respectively. Employing GPT-3.5 in Cantor also achieves an accuracy gain of 2.24% and 9.2%. In all of our experiments, we use only one MLLM (Gemini) to play the role of multiple experts, performing different sub-tasks with different requirements. Our contributions are the following:

- We propose an inspiring multimodal CoT framework named Cantor, which features a perceptual decision architecture that effectively integrates visual context and logical reasoning to solve visual reasoning tasks.
- We utilize the advanced cognitive abilities of an MLLM to act as multifaceted experts, obtaining higher-level information and significantly enhancing CoT generation.
- We demonstrate Cantor’s effectiveness on two challenging benchmarks, largely surpassing existing counterparts.

2. Related Work

2.1. Multimodal Large Language Models

Recent researches indicate that the development of Multimodal Large Language Models (MLLMs) [6, 10, 11, 33, 37, 39, 47, 48] is the result of combining the advanced reasoning capabilities of Large Language Models (LLMs) with the capabilities of Vision-Language models (VLMs). These models have achieved significant performance improvements in multimodal tasks by integrating visual and linguistic information. In particular, significant progress [13, 23, 36] has been made in connecting visual and text representations with contrastive visual and language models, but they encounter limitations when dealing with downstream tasks that require generating components or performing more refined reasoning on visual and language. To overcome these limitations, MLLM extends the reasoning and generation capabilities of LLM to the visual domain by directly inferring embedded visual features [1, 2, 7, 9, 24, 54]. In addition, MLLMs further improve performance through fine-tuning visual instructions [28].

These advances not only demonstrate the ability of MLLM to handle complex multimodal information but also provide new possibilities for achieving General Artificial Intelligence (AGI) with rich multimodal information. By integrating the text reasoning ability of LLM with the image understanding ability of visual language models, MLLM can achieve deep understanding and expression in multiple modalities, processing complex tasks such as image captioning and visual question answering. Open-source MLLMs such as LLaVA [28] demonstrate these capabilities, while closed-source models such as GPT4-V [34] and Gemini [40] have taken a greater step in capturing scene context, reasoning, and creativity. Although for specific tasks these closed-source models may not be directly com-

petent or fine-tuning. However, prompt learning can to some extent overcome these limitations. This paper is dedicated to exploring the technique of CoT [43] to enhance the ability of MLLMs to capture the complete context of complex visual scenes, thereby further strengthening their reasoning capabilities.

2.2. Tool-Augmented Language Models

In recent years, despite the impressive performance of Large Language Models (LLMs), they are not without their inherent limitations. These include challenges such as obtaining up-to-date information [21], the inability to employ specific tools [31, 38], and difficulties in executing complex reasoning processes [29, 30]. Meanwhile, researchers are increasingly interested in using external tools and modular methods to enhance LLM through prompting and in-context learning. These enhanced LLMs can utilize different external tools to provide LLMs with more functionality and gain more knowledge. Some works [5, 12, 17, 19] utilized prompts to generate complex programs that can be executed by computers, calling different tools to more effectively perform logical reasoning tasks. For example, PaLI-X-VPD [17] extracted the reasoning ability of LLM by generating multiple candidate programs, executing programs through external tools, and verifying their correctness. It transformed each correct program into a language description of reasoning steps to form a CoT. In addition, some works proposed benchmarks (such as API Bank [25], ToolQA [55], and MetaTool [18]) to evaluate the effectiveness of LLM tool use. This article mainly emphasizes enhancing the tool usage ability of MLLM.

2.3. Multi-modal CoT Reasoning

LLMs and MLLMs are becoming increasingly popular. Although their own abilities are becoming stronger, good prompt methods are still the key to fully unleashing their abilities. Chain-of-thought (CoT) is a method to improve LLM’s reasoning ability, and the core of CoT is to encourage LLM to clarify their reasoning in a human thinking way, specifically by adding logical thinking processes before obtaining answers. In the field of NLP, CoT has received extensive research [8, 15, 42, 51]. Jason Wei *et al.* [43] significantly improved LLM’s reasoning ability by simply adding problem-solving ideas directly to in-context examples. Subsequently, researchers mainly focused on how to automate the construction of CoT to reduce manual annotation and more complex structures such as Tree-of-Thought (ToT) [45] and Graph-of-Thought (GoT) [3, 22, 46].

Meanwhile, surprising progress has been made in multimodal CoT. MM-CoT [52] firstly proposed a two-stage reasoning framework by using text and image pairs as input, generating rationale first and then generating answers.

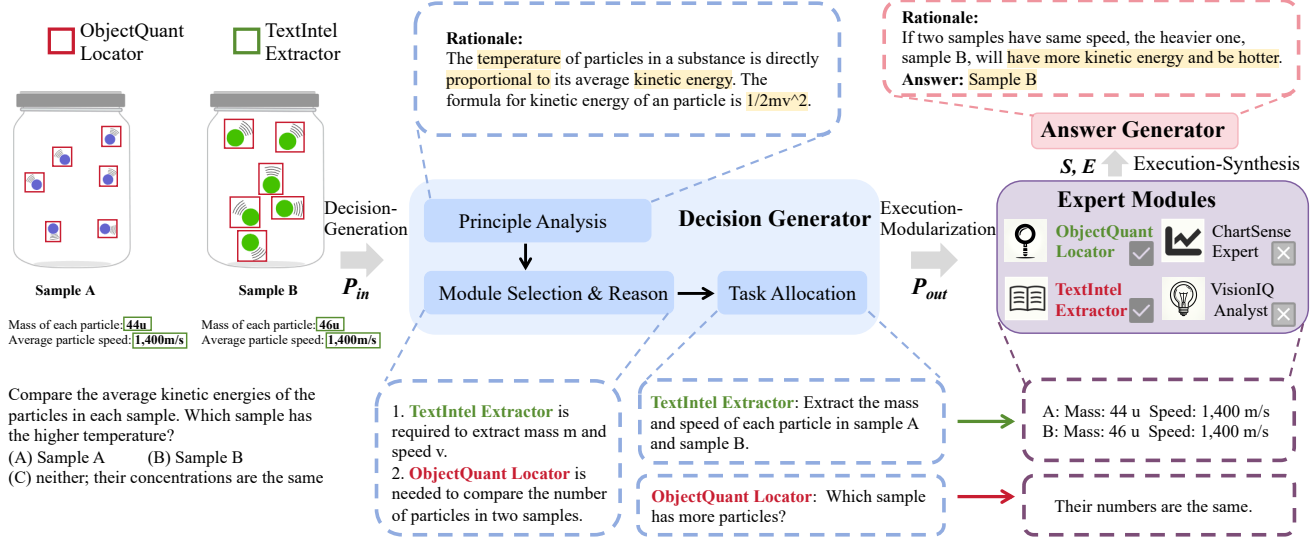


Figure 2. Overview of Cantor and a specific example. Cantor analyzes the image and problem through the Decision Generator, offering the principle analysis of the questions, and providing module selection & Reason, as well as specific task allocation. Subsequently, MLLM acts as various expert modules to execute sub-tasks. Finally, Cantor synthesizes and contemplates through the Answer Generator, providing the final answer.

Subsequent works [14, 14, 41, 53] are mostly based on this framework, focusing on designing special vision-language feature fusion mechanisms to enhance multimodal information interaction. However, these CoT prompting methods need to fine-tune on ground truth of natural language reasoning, which requires both annotation and computation costly. Based on this issue, researchers have proposed other CoT methods that do not require manual annotation and training. On the one hand, they fully tap into textual information. For example, DD-CoT [53] further refined the process of generating the CoT. Without introducing visual information, it used LLM to break down the problem into multiple related sub-questions and then answer each sub-question one by one to form the CoT. On the other hand, researchers are committed to enhancing visual information through various means. For example, CoCoT [49] captured image characteristics by comparing the similarities and differences between images, while CCoT [32] obtained scene maps by disassembling the targets and attributes in the images to assist in rationale generation. The key difference between our method and these methods is that when mining text information, we introduce visual information in advance to make decisions more reasonable and factual. In addition, we enhance visual information more comprehensively by calling multiple experts. Last, Cantor is also a method that does not require training or manual annotation, so it has strong universality and convenience. This paper emphasizes enhancing the expert usage capability of MLLM. Considering that MLLM has multimodal universal capabilities, it is naturally suitable to serve as various ex-

perts. Therefore, this paper will endow MLLM with various identities and explore its expert-playing abilities.

3. Method

To address the limitations of multimodal CoT in solving visual reasoning tasks, we propose Cantor, which introduces visual information to make correct decisions and uses a single MLLM to act as multiple experts to adapt to a wide range of problems. We describe the framework of Cantor (Section 3.1). Then, we provide a detailed introduction to our two-step approach: the first is Decision-Generation (Section 3.2), and the second is Execution (Section 3.3).

3.1. Preliminaries

Cantor consists of two stages: Decision-Generation and Execution, as shown in Fig. 2. During the Decision-Generation stage in Cantor, Cantor’s input consists of $X = \{I, T, P_{in}\}$, where I denotes the visual input (image or a caption), T signifies the text input, which represents the concatenation of the problem statement and its context, and P_{in} represents the prompt for generating decisions. Formally, given an input query X , a decision P is generated as follows: $P_{out} = F(X)$, where F denotes the decision generator (an LLM or MLLM). Specially, $P_{out} = \{R, O, S_t\}$, where R denotes Principle Analysis, O denotes Module Selection & Reason, and S_t denotes the tasks assigned to expert modules. For specific examples, please refer to the blue section in the middle of Fig. 2.

In the execution-modularization stage, multiple sub-

tasks $S_t = \{st_1, st_2 \dots st_n\}$ derived from the decision P_{out} and image I are jointly sent to the corresponding expert module to obtain the sub-answers $S_a = \{sa_1, sa_2, \dots, sa_n\}$. The process is as follows: $S_a = G(S_t, I)$, where G denotes various experts (an MLLM). This process corresponds to the Execution-Modularization stage in the purple section at the bottom right of Fig. 2. Then in Execution-Synthesis stage, we concatenate the sub-tasks and sub-answers to form supplementary information $S = \{S_t, S_a\}$, and design an answer generation prompt E . Finally, feed the updated input $X' = \{I, T, S, E\}$ and infer the final answer $A = F(X')$, where F denotes the answer generator (an LLM or MLLM), as shown in the upper right corner of Fig. 2.

3.2. Step 1: Decision-Generation

Our first step is to generate decision P_{out} which considers and deploys the problem. Please note that we are studying unsupervised visual reasoning tasks, which involve having the model generate corresponding decisions for the problem without ground truth [44, 49]. Additionally, for standardization and accuracy, we adopt a few-shot setting in prompt to provide a decision generation prompt P_{in} for the model, which includes the requirements for decision generation, the characteristics of callable modules, and several manually written decision examples.

Let’s provide a detailed introduction to the Decision-Generation process of Cantor and the specific components of the prompt P_{in} :

1. Acting as Decision Generator. We prompt the LLM or MLLM with “You are an advanced question-answering agent required with four specialized modules to aid in the analysis and responding to queries about images” enabling it to function as a decision generator in Cantor.

2. Expert Modules Unveiled. As shown in the Expert Modules of Fig. 2. We provide detailed information on the characteristics of each expert module for Cantor, with the aim to allocate tasks to each expert module based on the principle of addressing the problem during the Decision-Generation phase, as follows: **TextIntel Extract:** This module extracts and converts text within images into editable text format. It’s particularly useful for images containing a mix of text and graphic elements. **ObjectQuant Locator:** This module identifies and locates objects within an image. It’s advanced at comparing quantities and recognizing spatial relationships. **VisionIQ Analyst:** This module processes and interprets visual data, enabling you to ask any queries related to the image’s content. **ChartSense Expert:** This module specializes in analyzing and interpreting information from charts and graphs. It can extract data points, understand trends, and identify key components such as titles, axes, labels, and legends within a chart.

3. Principle Analysis and Module Selection & Rea-

son. We prompt Cantor “Provide a rationale for your approach to answering the question, explaining how you will use the information from the image and the modules to form a comprehensive answer”, performing an overall assessment and modular analysis of the question.

4. Task Allocation. We prompt “Assign specific tasks to each module as needed, based on their capabilities, to gather additional information essential for answering the question accurately.”, requiring Cantor to select the necessary modules and assign their corresponding specific tasks.

5. Contextual Insights and Practical Applications. We introduce some in-context examples to enhance Cantor’s comprehension of our prompts, ensuring its responses adhere to the desired format. Detailed instances are provided in the supplementary materials for further reference. Then, we input the particular problem that needs addressing, along with its contextual details, enabling Cantor to formulate nuanced decisions. The blue part on the left half of Fig. 2 shows a specific example of decision generation.

The above five parts are combined to form the final decision generation prompt P_{in} . Subsequently, P_{in} together with visual input I and text input T , constitutes the complete input for the first stage of Cantor, prompting Cantor to deliver a deliberate decision P_{out} .

The decision generation method represents a core novel contribution of our work. Initially, the LLM or MLLM is employed as a decision generator, serving as the brain. Next, a suite of specialized expert modules is integrated, augmenting the decision generating with diverse capabilities analogous to the limbs. This integration ensures that decision-generating is both comprehensive and granular, leveraging the strengths of each module. Thereafter, the decision generator tailors tasks for selected expert modules based on insights gained from principle analyses. This dynamic task allocation enhances Cantor’s efficiency and effectiveness. Ultimately, the introduction of in-context examples enables the MLLM to learn and reference, thereby further improving the accuracy and adaptability of decision generation. Notably, we introduce visual context in advance during the Decision-Generation stage, rather than the Execution stage, effectively alleviating determining hallucinations.

3.3. Step 2: Execution

In Cantor, the execution stage can be divided into two stages, Execute-Modularization and Execute-Synthesis. The former completes the sub-tasks assigned during the Decision-Generation stage by calling various expert modules and providing supplementary information. The latter summarizes various supplementary information from the execute-modularization stage and generates the final answer through rational and detailed thinking.

Execute-Modularization. We call the expert mod-

ule to execute the various sub-tasks assigned during the Decision-Generation stage. Specially, we first extract sub-tasks $S_t = \{st_1, st_2 \dots st_n\}$ from P_{out} . Next, we find the expert module corresponding to the sub-task st_i in sequence, and input the sub-task st_i as the prompt into the expert, such as “ObjectQuant Locator: Which sample has more particles?”. Subsequently, we obtain the sub-task answer sa_i , such as “Their numbers are the same”, as shown in the lower right part of Fig. 2.

Symbolically, we input the experts played by MLLM, sub-task st_i , and image I , and MLLM provides the execution results of the sub-task. The process is as follows: $sa_i = G(I, st_i)$, where $G(\cdot)$ represents MLLM acting as experts, and sa_i represents the sub-task’s answer. When executing sub-tasks, we only use one MLLM to act as different expert modules. This not only simplifies the pipeline of the method but also aims to fully utilize the advanced cognitive abilities of MLLM.

Execute-Synthesis. We concatenate and summarize the obtained sub-tasks and sub-tasks answers to obtain supplementary information S for auxiliary reasoning, as follows: $S = \{[st_1, sa_1] \cdot [st_2, sa_2] \cdot \dots \cdot [st_n, sa_n]\}$. Notably, in the answer generation stage, we introduce the answer generation prompt E , which includes the prompt and the formatting requirement for generating answers, as follows: “You are a knowledgeable and skilled information integration science expert. Please gradually think and answer the questions based on the given questions, options, and supplementary information. Please note that we not only need answers but more importantly, we need rationales for obtaining answers. Please combine your knowledge and supplementary information to obtain reasoning and answers. Please prioritize using your knowledge to answer questions. If unable to answer, maintain critical thinking and select effective information to assist you in selecting the most correct option as the answer. Furthermore, please do not rely solely on supplementary information, as the provided supplementary information may not always be effective.”

This includes three key points. Firstly, we use prompts to have Cantor play the role of an answer generator who is knowledgeable and skilled at integrating information. This not only ensures its professionalism and ability to make basic judgments on questions but also ensures that it can better integrate information obtained during the Execute-Modularization stage. Secondly, to increase interpretability, demonstrate the thinking process of Cantor, and improve its thinking ability, we require Cantor to answer the basic principles first, and then generate the corresponding options, as shown in the pink box in Fig. 2. Finally, we request that Cantor remain rational and critical, ensuring it does not solely rely on the information obtained from the Execute-Modularization stage. This approach promotes a more balanced and comprehensive execute-synthesis process.

4. Experiments

In this section, we evaluate the proposed Cantor on two visual reasoning datasets: ScienceQA [29] and MathVista [30]. The experimental results show that Cantor outperforms existing baselines in these tasks. Additionally, we analyze the importance of visual information in visual reasoning tasks. Finally, we conduct a detailed analysis of Cantor’s key components.

4.1. Datasets

We evaluate our method on two visual reasoning task benchmarks.

ScienceQA [29]: It is the first multimodal scientific question-and-answer dataset annotated with detailed explanations. The problems with datasets are systematically divided into three main scientific disciplines: natural sciences (NAT), social sciences (SOC), and language sciences (LAN). We only use the ScienceQA test set, which contains 4241 questions and answers, of which 2,017 samples are attached with images.

MathVista [30]: It is a dataset that combines the challenges of various mathematical and visual tasks. It requires high levels of model granularity, deep visual understanding, and combinatorial reasoning ability, making it a challenging dataset for current basic models. In the experiment, we used Mathvista testmini, which includes 1000 text and image pairs for Q&A.

4.2. Models

We use two models to evaluate our method, GPT-3.5 and Gemini Pro 1.0, by calling their official API. Firstly, we use GPT-3.5 to evaluate the impact of introducing high-level perceptual information on LLM inference ability and explore the linkage ability between LLM and MLLM. Secondly, we use Gemini Pro 1.0, an advanced MLLM. We desire to fully tap into the multimodal ability of MLLM and improve its reasoning ability.

4.3. Implementation Details

We implement two versions of Cantor based on GPT-3.5 and Gemini. Cantor(GPT-3.5) uses both GPT-3.5 as the Decision Generator and Answer Generator during the Decision-Generation and Execute-Synthesis stage. Differently, Cantor(Gemini) uses Gemini in these two stages. For the Execute-Modularization stage, due to the need for multimodality, we use Gemini as the MLLM in both versions, playing various roles as experts. For the captions required for Cantor(GPT-3.5) in the Decision-Generation stage, we generated them through Gemini Pro 1.0, with the prompt “Please provide the detailed title of this image as much as possible”. In terms of models’ prompts, although the two models have different preferences for prompts, we

| Methods | Supervised | IMG | NAT | SOC | LAN | TXT | NO | G1-6 | G7-12 | Avg |
|---------------------------|------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Random Chance | ✗ | 40.08 | 40.28 | 46.13 | 29.25 | 47.45 | 33.66 | 39.35 | 40.67 | 39.83 |
| Human Average [29] | ✗ | 87.50 | 90.23 | 84.97 | 87.48 | 89.60 | 88.10 | 91.59 | 82.42 | 88.40 |
| UnifiedQA [20] | ✓ | 61.38 | 68.16 | 69.18 | 74.91 | 63.78 | 77.84 | 72.98 | 65.00 | 70.12 |
| UnifiedQA (CoT) [20] | ✓ | 66.53 | 71.00 | 76.04 | 78.91 | 66.42 | 81.81 | 77.06 | 68.82 | 74.11 |
| Multimodal-CoT [52] | ✓ | 82.90 | 87.52 | 77.17 | 85.82 | 87.88 | 86.83 | 84.65 | 85.37 | 84.91 |
| LLaMA-Adapter [50] | ✓ | 80.32 | 84.37 | 88.30 | 84.36 | 83.72 | 86.90 | 85.83 | 84.05 | 85.19 |
| LLaVa [28] | ✓ | 88.00 | 90.36 | 95.95 | 88.00 | 89.49 | 90.66 | 90.93 | 90.90 | 90.92 |
| LLaVA (GPT-4) [28] | ✓ | 88.99 | 91.56 | 96.74 | 91.09 | 90.62 | 93.52 | 92.73 | 92.16 | 92.53 |
| LLaMA-SciTune (CTOM) [16] | ✓ | 86.67 | 89.30 | 95.61 | 87.00 | 93.08 | 91.75 | 84.37 | 91.30 | 90.03 |
| GPT-3 (zero-shot) [4] | ✗ | 65.74 | 75.04 | 66.59 | 78.00 | 74.24 | 79.58 | 76.36 | 69.87 | 74.04 |
| GPT-3.5 (CoT) (AE) [35] | ✗ | 66.09 | 76.60 | 65.92 | 77.55 | 75.51 | 79.58 | 78.49 | 67.63 | 74.61 |
| GPT-3.5 (CoT) (ALE) [35] | ✗ | 67.43 | 75.44 | 70.87 | 78.09 | 74.68 | 79.93 | 78.23 | 69.68 | 75.17 |
| GPT-3.5 CoT [33] | ✗ | 67.92 | 78.82 | 70.98 | <u>83.18</u> | 77.37 | 86.13 | 80.72 | 74.03 | 78.31 |
| QVix(GPT-3.5) [44] | ✗ | 55.00 | - | - | - | - | - | - | - | - |
| Chameleon (GPT-3.5) [31] | ✗ | 70.80 | 81.62 | 70.64 | 84.00 | 79.77 | <u>86.62</u> | 81.86 | <u>76.53</u> | 79.93 |
| DD-CoT(GPT-3) [53] | ✗ | 69.96 | 78.60 | 73.90 | 80.45 | 77.27 | 82.93 | 80.65 | 73.50 | 78.09 |
| DD-CoT(GPT3.5) [53] | ✗ | <u>72.53</u> | 80.15 | <u>76.72</u> | 82.82 | <u>78.89</u> | 85.02 | <u>82.86</u> | 75.21 | <u>80.15</u> |
| Cantor(GPT-3.5) | ✗ | 77.54 | <u>80.37</u> | 85.49 | 84.00 | 77.27 | 86.83 | 85.61 | 76.60 | 82.39 |
| Gemini | ✗ | 76.85 | 79.13 | 85.26 | 80.82 | 76.93 | 83.83 | 83.81 | 75.54 | 80.85 |
| Cantor(Gemini) | ✗ | 82.40 | 84.24 | 87.85 | 84.09 | 82.11 | 86.97 | 88.18 | 79.17 | 84.96 |

Table 1. Accuracy scores (%) on ScienceQA [29], where **bold** entries indicate the best results, underlines indicate the second-best. We compare the performance of our system with various baseline models including supervised models and unsupervised models. Question classes: NAT = natural science, SOC = social science, LAN = language science, TXT = text context, IMG = image context, NO = no context, G1-6 = grades 1-6, G7-12 = grades 7-12.

use the same prompt for the sake of method universality in Decision-Genetation stage and Execute-Synthesis stage. The prompt in Execute-Modularization stage is generated by the Cantor itself. For different datasets’ prompts, we design different in-context examples based on their characteristics, and the rest of the prompts are the same.

4.4. Main Results

ScienceQA. Tab. 1 shows the results of existing baselines compared to our method Cantor on ScienceQA. Using GPT-3.5 as the base LLM to decision and answer, Cantor achieves an accuracy of 82.39%, which is an improvement of 4.08% over the chain-of-thought (CoT) prompted GPT-3.5 [33]. Furthermore, with Gemini as the decision generator and answer generator, Cantor reaches an accuracy of 84.96%, significantly surpassing all training-free methods, and even outperforming fine-tuned methods like UnifiedQA (CoT) [52] and MM-CoT [52]. This not only demonstrates the generality of Cantor but also shows that Cantor starts with perception-based information for making better decisions. Moreover, by invoking various expert modules, it can introduce richer contextual information to both LLMs and MLLMs, aiding in problem-solving.

Particularly noteworthy is that Cantor advances in the multimodal domain. As shown in Tab. 2, we further present

| Method | Subject | | | Grade | | Average |
|---------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | NAT | SOC | LAN | G1-6 | G7-12 | |
| LLaVA | 37.0 | 61.5 | 33.3 | 52.3 | 30.5 | 46.2 |
| MiniGPT | 45.2 | 51.5 | 38.1 | 50.6 | 39.1 | 47.4 |
| InstructBLIP | 43.9 | 58.1 | 47.6 | 53.1 | 39.4 | 49.3 |
| QVix (GPT-3.5) | 48.0 | 67.1 | 38.1 | 60.6 | 40.5 | 55.0 |
| Qwen-VL-Chat | - | - | - | - | - | 68.85 |
| mPLUG-Ow12 | - | - | - | - | - | 68.75 |
| Chameleon (GPT-3.5) | - | - | - | - | - | 70.8 |
| SPHINX-2k | - | - | - | - | - | 70.6 |
| LLaVA1.5 | - | - | - | - | - | 71.6 |
| GPT-3.5 (+Caption) | 70.14 | 62.43 | 68.18 | 78.59 | 52.32 | 67.18 |
| Cantor (GPT-3.5) | 73.45 | 83.38 | 88.64 | 84.31 | 66.55 | 77.54 |
| Gemini | 71.55 | 84.29 | 93.18 | 80.90 | 67.01 | 76.85 |
| Cantor (Gemini) | 79.49 | 86.39 | 93.18 | 86.98 | 71.26 | 82.40 |

Table 2. Accuracy scores (%) on ScienceQA for the IMG class, which includes image context.

the accuracy of various methods on ScienceQA for the IMG class, which includes image context. It can be seen that Cantor based on GPT-3.5 significantly surpasses the base-

| Model | Input | FQA | GPS | MWP | TQA | VQA | ALG | ARI | GEO | LOG | NUM | SCI | STA | ALL |
|---|--|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| <i>Heuristics baselines</i> | | | | | | | | | | | | | | |
| Random chance | - | 18.2 | 21.6 | 3.8 | 19.6 | 26.3 | 21.7 | 14.7 | 20.1 | 13.5 | 8.3 | 17.2 | 16.3 | 17.9 |
| Frequent guess | - | 22.7 | 34.1 | 20.4 | 31.0 | 24.6 | 33.1 | 18.7 | 31.4 | 24.3 | 19.4 | 32.0 | 20.9 | 26.3 |
| <i>Large Language Models (LLMs)</i> | | | | | | | | | | | | | | |
| Zero-shot GPT-3.5 | <i>Q</i> only | 21.9 | 26.9 | 9.1 | 38.6 | 23.5 | 27.7 | 15.9 | 25.7 | 21.6 | 9.9 | 41.5 | 20.5 | 23.5 |
| Zero-shot GPT-4 | <i>Q</i> only | 22.3 | 37.0 | 7.0 | 39.2 | 27.4 | 33.6 | 17.4 | 35.6 | 16.2 | 9.2 | 45.8 | 19.5 | 26.1 |
| Zero-shot Claude-2 | <i>Q</i> only | 21.9 | 34.1 | 13.4 | 36.1 | 29.1 | 32.8 | 20.4 | 33.3 | 13.5 | 12.1 | 36.4 | 20.5 | 26.4 |
| <i>Augmented Large Language Models (Augmented-LLMs)</i> | | | | | | | | | | | | | | |
| 2-shot CoT GPT-3.5 | <i>Q, I_c, I_t</i> | 27.5 | 29.3 | <u>36.0</u> | 49.4 | 29.1 | 31.0 | <u>32.9</u> | 31.0 | <u>16.2</u> | 17.4 | 50.8 | 37.2 | 33.2 |
| 2-shot CoT GPT-4 | <i>Q, I_c, I_t</i> | 27.9 | 31.7 | 31.2 | <u>51.9</u> | 28.5 | 33.5 | 30.9 | 32.2 | 13.5 | 12.5 | 58.2 | <u>37.9</u> | 33.2 |
| 2-shot PoT GPT-3.5 | <i>Q, I_c, I_t</i> | 24.5 | 26.4 | 23.7 | 33.5 | 27.9 | 27.8 | 26.1 | 28.0 | 18.9 | 13.2 | 33.6 | 29.9 | 26.8 |
| 2-shot PoT GPT-4 | <i>Q, I_c, I_t</i> | <u>30.1</u> | 39.4 | 30.6 | 39.9 | <u>31.3</u> | 37.4 | 31.7 | 41.0 | 18.9 | <u>20.1</u> | 44.3 | <u>37.9</u> | <u>33.9</u> |
| GPT-3.5 | <i>Q, I_c</i> | 26.0 | 31.7 | 35.5 | 48.1 | 30.2 | 32.4 | 32.3 | 33.0 | <u>16.2</u> | 17.4 | 54.9 | 36.2 | 33.2 |
| Cantor (GPT-3.5) | <i>Q, I_c</i> | 45.7 | <u>31.8</u> | 40.9 | 55.1 | 44.1 | <u>34.5</u> | 42.2 | <u>33.9</u> | 13.5 | 36.1 | <u>55.0</u> | 55.5 | 43.1 |
| <i>Multimodal Large Language Models (MLLMs)</i> | | | | | | | | | | | | | | |
| IDEFICS-9B-Instruct | <i>Q, I</i> | 21.6 | 21.1 | 6.5 | 25.9 | 24.0 | 22.1 | 15.0 | 19.8 | <u>18.9</u> | 9.9 | 24.6 | 18.1 | 19.8 |
| mPLUG-Owl-LLaMA-7B | <i>Q, I</i> | 22.7 | 23.6 | 10.2 | 27.2 | 27.9 | 23.6 | 19.2 | 23.9 | 13.5 | 12.7 | 26.3 | 21.4 | 22.2 |
| miniGPT4-LLaMA-2-7B | <i>Q, I</i> | 18.6 | 26.0 | 13.4 | 30.4 | 30.2 | 28.1 | 21.0 | 24.7 | 16.2 | 16.7 | 25.4 | 17.9 | 23.1 |
| LLaMA-Adapter-V2-7B | <i>Q, I</i> | 21.2 | 25.5 | 11.3 | 32.3 | 31.8 | 26.3 | 20.4 | 24.3 | 24.3 | 13.9 | 29.5 | 18.3 | 23.9 |
| LLaVAR | <i>Q, I</i> | 21.9 | 25.0 | 16.7 | 34.8 | 30.7 | 24.2 | 22.1 | 23.0 | 13.5 | 15.3 | 42.6 | 21.9 | 25.2 |
| InstructBLIP-Vicuna-7B | <i>Q, I</i> | 23.1 | 20.7 | 18.3 | 32.3 | 35.2 | 21.8 | 27.1 | 20.7 | <u>18.9</u> | 20.4 | 33.0 | 23.1 | 25.3 |
| LLaVA-LLaMA-2-13B | <i>Q, I</i> | 26.8 | 29.3 | 16.1 | 32.3 | 26.3 | 27.3 | 20.1 | 28.8 | 24.3 | 18.3 | 37.3 | 25.1 | 26.1 |
| Multimodal Bard | <i>Q, I</i> | 26.0 | 47.1 | 29.6 | 48.7 | 26.8 | 46.5 | 28.6 | 47.8 | 13.5 | 14.9 | <u>47.5</u> | 33.0 | 34.8 |
| Gemini | <i>Q, I</i> | <u>37.1</u> | 29.3 | <u>38.1</u> | 57.5 | <u>36.3</u> | 36.0 | <u>35.7</u> | 31.4 | 24.3 | <u>25.7</u> | 50.0 | 41.9 | <u>38.8</u> |
| Cantor (Gemini) | <i>Q, I</i> | 50.2 | <u>39.4</u> | 39.8 | <u>49.4</u> | 43.8 | <u>42.0</u> | 41.5 | <u>41.4</u> | 10.8 | 30.8 | 46.7 | 59.5 | 44.7 |

Table 3. Accuracy scores (%) on the *testmini* subset of MathVista, where **bold** entries indicate the best results, underlines indicate the second-best. Input: *Q*: question, *I*: image, *I_c*: image caption, *I_t*: OCR text detected in the image. Task types: FQA: figure question answering, GPS: geometry problem solving, MWP: math word problem, TQA: textbook question answering, VQA: visual question answering. Mathematical reasoning types: ALG: algebraic reasoning, ARI: arithmetic reasoning, GEO: geometry reasoning, LOG: logical reasoning, NUM: numeric commonsense, SCI: scientific reasoning, STA: statistical reasoning. ALL: overall accuracy. The performance results in the table come from [30].

line in various problems, and even surpasses well-known MLLMs such as SPHINX [26] and LLaVA-1.5 [27]. This indicates that clear perceptual decisions can trigger the reasoning ability of language models toward dense image information. At the same time, the experiment on Gemini also shows that we further stimulate the visual reasoning ability of MLLM.

MathVista. MathVista [30] is a challenging dataset that integrating a variety of mathematical reasoning tasks with visual tasks. Tab. 3 compares different method performances. We also conduct experiments using GPT-3.5 and Gemini as baselines. From general visual question answering to professional math word problems, Cantor has greatly surpassed the baseline in almost all types of prob-

| Analysis | ScienceQA | MathVista |
|-----------------------|--------------|--------------|
| No Visual Information | 65.69 | 25.70 |
| + Rough Caption | 63.21 | 25.10 |
| + Detailed Caption | 74.37 | 33.20 |
| + Image | 78.85 | 38.00 |

Table 4. The impact of different levels of visual information on model’s performance.

lems. This indicates that correct decision and modular experts can stimulate their fine-grained, in-depth visual understanding and combinatorial reasoning abilities. It is worth noting that Cantor (GPT-3.5) even surpasses GPT-4 based

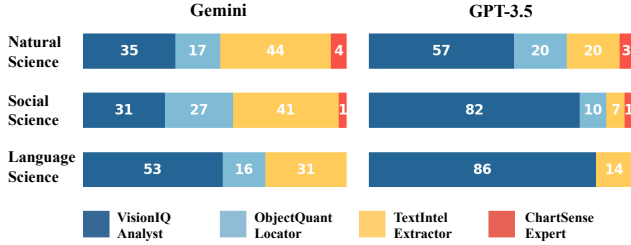


Figure 3. Proportions of Cantor’s invocation of expert modules across three types of questions on ScienceQA.

on CoT and PoT.

4.5. Quantitative Analysis

Analysis Visual Cues for Decision Generation. We conduct a detailed analysis of the impact of visual information on Gemini’s decision generation on ScienceQA and MathVista, with the prompt “think step by step”. The results are shown in Tab. 4. When we do not input any form of visual information (including images and captions) in the experiment, only the text of the question is input. It can be seen that even without any visual information, MLLMs like Gemini still possess strong logical reasoning ability in pure language modal, demonstrating its superiority as a decision generator. Then we step by step explore the impact of incorporating visual information on Gemini. Firstly, we add rough captions, such as “A photo of a black and white cat.” Gemini’s performances unexpectedly decline on both datasets. This indicates that overly simplistic captions not only fail to promote MLLM, but can even mislead them into making incorrect decisions. Next, we enrich the description of captions to fully reproduce the image scene as much as possible. It can be seen that with the addition of detailed captions, Gemini’s performance has significantly improved compared to those without visual information or rough captions. This indicates that visual information is indispensable for complex visual reasoning tasks. Finally, we replace captions with images, and it can be seen that Gemini’s performance increased by 4.48% and 4.8% on both datasets, achieving the best performance at the same time. This is also in line with intuition, as the generation of captions is uncontrollable and may not necessarily contain key information for solving problems, but images themselves must have complete information. Therefore, in complex visual reasoning tasks, using images instead of captions to obtain visual information is a better solution for MLLM.

Expert Module Use Planning. The proportion of Cantor calling various expert modules on ScienceQA is shown in Fig. 3. We find that GPT-3.5 and Gemini exhibit different decision-generating behaviors. GPT-3.5 has a strong preference for using Object Quant Locator, with usage rates exceeding 80% in both Social Science and Language Science

| Module | Enable Only | Disable Only |
|---------------------|--------------|--------------|
| TextIntel Extractor | 80.91(+4.06) | 80.86(-1.54) |
| ObjectQuant Locator | 80.27(+3.42) | 81.01(-1.39) |
| VisionIQ Analyst | 80.22(+3.37) | 81.51(-0.89) |
| ChartSense Expert | 79.13(+2.28) | 81.71(-0.69) |
| Gemini / Cantor | 76.85 | 82.40 |

Table 5. Performance increase with enabled modules and performance drop with disabled modules on ScienceQA, where “Enable Only” only just this module is on, others off. “Disable Only” means just this module is off, others on. In the last line, “Gemini/Cantor” denotes the original Gemini baseline and the fully implemented version of Cantor.

subjects, far exceeding other expert modules. We speculate that this is because GPT-3.5 is heavily influenced by in-context examples. On the other hand, Gemini is relatively balanced in expert module calls and does not exhibit any particular preferences. In addition, the usage ratio of both modules for ChartSense Expert is very low, especially for the Language Science subject where the number of calls is 0. This is because the proportion of questions related to table content is very small in ScienceQA, and there is even no question about table content in Language Science. This demonstrates the rationality of the decisions made by the two models. For different types of problems, the Language Science subject focuses more on the language meaning behind the image rather than being limited to the combination of target numbers or positions. Therefore, the two models call VisionIQ Analyst more frequently, reducing the use of ObjectQuant Locator.

Ablation Study with Modules. We use Gemini as the MLLM to investigate the impact of enabling and disabling expert modules on the performance of ScienceQA. The results are shown in Tab. 5. The results show that the use of each expert module results in a gain (maximum 4.06%, minimum 2.28%), indicating that all expert modules play a crucial role. The TextIntel Extractor is the most important among all modules, with the most significant gains and decreases in performance. At the same time, we can also find that enabling a module has a greater impact on model performance than disabling it. We believe that the effective high-level information obtained by an expert module (MLLM) is more generalized, compared with lower-level visual-information (such as coordinates, color, attributes, etc.). This higher-level information assists in the execution of other module tasks. In our method, even if a module is disabled, MLLM playing the role of other experts can to some extent compensate for the lack of that module, as they are not operating in isolation. We have also added some results in the supplementary material to support this view.

5. Conclusion

In this paper, we introduce an inspiring multimodal chain-of-thought framework named Cantor, designed to enhance the determining capabilities of MLLMs. By delving into the pivotal role of visual information in the decision-generating process, this paper highlights the importance of integrating visual cues at the decision stage, effectively mitigating the hallucination issues that may arise in LLMs. The novelty of the Cantor framework also lies in its ability to enable an MLLM to emulate the roles of domain-specific experts, acquiring high-level information, and thereby facilitating more rational and in-depth reasoning processes. Demonstrated on the challenging benchmarks of ScienceQA and MathVista involving complex visual reasoning tasks, Cantor has shown remarkable adaptability and efficacy, proving its strong potential in addressing real-world reasoning problems across various domains.

References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022. 3
- [2] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023. 3
- [3] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 17682–17690, 2024. 3
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020. 7
- [5] Wenhui Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv preprint arXiv:2211.12588*, 2022. 3
- [6] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240): 1–113, 2023. 3
- [7] 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. 3
- [8] Shizhe Diao, Pengcheng Wang, Yong Lin, and Tong Zhang. Active prompting with chain-of-thought for large language models. *arXiv preprint arXiv:2302.12246*, 2023. 3
- [9] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023. 3
- [10] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023. 3
- [11] Chaoyou Fu, Renrui Zhang, Zihan Wang, Yubo Huang, Zhengye Zhang, Longtian Qiu, Gaoxiang Ye, Yunhang Shen, Mengdan Zhang, Peixian Chen, Sirui Zhao, Shaohui Lin, Deqiang Jiang, Di Yin, Peng Gao, Ke Li, Hongsheng Li, and Xing Sun. A challenger to gpt-4v? early explorations of gemini in visual expertise. *arXiv preprint arXiv:2312.12436*, 2023. 3
- [12] Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine Learning*, pages 10764–10799. PMLR, 2023. 3
- [13] Shashank Goel, Hritik Bansal, Sumit Bhatia, Ryan Rossi, Vishwa Vinay, and Aditya Grover. Cyclip: Cyclic contrastive language-image pretraining. *Advances in Neural Information Processing Systems*, 35:6704–6719, 2022. 3
- [14] Liqi He, Zuchao Li, Xiantao Cai, and Ping Wang. Multimodal latent space learning for chain-of-thought reasoning in language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 18180–18187, 2024. 4
- [15] Namgyu Ho, Laura Schmid, and Se-Young Yun. Large language models are reasoning teachers. *arXiv preprint arXiv:2212.10071*, 2022. 3
- [16] Sameera Horawalavithana, Sai Munikoti, Ian Stewart, and Henry Kvinge. Scitune: Aligning large language models with scientific multimodal instructions. *arXiv preprint arXiv:2307.01139*, 2023. 7
- [17] Yushi Hu, Otilia Stretcu, Chun-Ta Lu, Krishnamurthy Viswanathan, Kenji Hata, Enming Luo, Ranjay Krishna, and Ariel Fuxman. Visual program distillation: Distilling tools and programmatic reasoning into vision-language models. *arXiv preprint arXiv:2312.03052*, 2023. 2, 3
- [18] Yue Huang, Jiawen Shi, Yuan Li, Chenrui Fan, Siyuan Wu, Qihui Zhang, Yixin Liu, Pan Zhou, Yao Wan, Neil Zhenqiang Gong, et al. Metatool benchmark for large language models: Deciding whether to use tools and which to use. *arXiv preprint arXiv:2310.03128*, 2023. 3
- [19] Shima Imani, Liang Du, and Harsh Shrivastava. Mathprompter: Mathematical reasoning using large language models. *arXiv preprint arXiv:2303.05398*, 2023. 3
- [20] Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and Hannaneh Hajishirzi.

- UnifiedQA: Crossing format boundaries with a single QA system. In *Findings of the Association for Computational Linguistics (EMNLP)*, pages 1896–1907, 2020. 7
- [21] Mojtaba Komeili, Kurt Shuster, and Jason Weston. Internet-augmented dialogue generation. *arXiv preprint arXiv:2107.07566*, 2021. 3
- [22] Bin Lei, Chunhua Liao, Caiwen Ding, et al. Boosting logical reasoning in large language models through a new framework: The graph of thought. *arXiv preprint arXiv:2308.08614*, 2023. 3
- [23] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pages 12888–12900. PMLR, 2022. 3
- [24] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 3
- [25] Minghao Li, Feifan Song, Bowen Yu, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. Api-bank: A benchmark for tool-augmented llms. *arXiv preprint arXiv:2304.08244*, 2023. 3
- [26] Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi Shao, Keqin Chen, et al. Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models. *arXiv preprint arXiv:2311.07575*, 2023. 8
- [27] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023. 8
- [28] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 3, 7
- [29] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*, 2022. 1, 3, 6, 7
- [30] 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. 1, 3, 6, 8
- [31] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 3, 7
- [32] Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. Compositional chain-of-thought prompting for large multimodal models. *arXiv preprint arXiv:2311.17076*, 2023. 2, 4
- [33] OpenAI. Chatgpt, 2022. 3, 7
- [34] OpenAI. GPT-4 technical report. *ArXiv*, abs/2303.08774, 2023. 3
- [35] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35: 27730–27744, 2022. 7
- [36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 3
- [37] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020. 3
- [38] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 3
- [39] Yi Tay, Mostafa Dehghani, Vinh Q Tran, Xavier Garcia, Jason Wei, Xuezhi Wang, Hyung Won Chung, Siamak Shakeri, Dara Bahri, Tal Schuster, et al. U12: Unifying language learning paradigms. *arXiv preprint arXiv:2205.05131*, 2022. 3
- [40] 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. 3
- [41] Lei Wang, Yi Hu, Jiabang He, Xing Xu, Ning Liu, Hui Liu, and Heng Tao Shen. T-sciq: Teaching multimodal chain-of-thought reasoning via large language model signals for science question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 19162–19170, 2024. 2, 4
- [42] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022. 3
- [43] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022. 3
- [44] Kaiwen Yang, Tao Shen, Xinmei Tian, Xiubo Geng, Chongyang Tao, Dacheng Tao, and Tianyi Zhou. Good questions help zero-shot image reasoning. *arXiv preprint arXiv:2312.01598*, 2023. 2, 5, 7
- [45] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [46] Yao Yao, Zuchao Li, and Hai Zhao. Beyond chain-of-thought, effective graph-of-thought reasoning in large language models. *arXiv preprint arXiv:2305.16582*, 2023. 3

- [47] 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. 3
- [48] Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. Woodpecker: Hallucination correction for multimodal large language models. *arXiv preprint arXiv:2310.16045*, 2023. 3
- [49] Daoan Zhang, Junming Yang, Hanjia Lyu, Zijian Jin, Yuan Yao, Mingkai Chen, and Jiebo Luo. Cocot: Contrastive chain-of-thought prompting for large multimodal models with multiple image inputs. *arXiv preprint arXiv:2401.02582*, 2024. 2, 4, 5
- [50] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Qiao Yu. LLaMA-Adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023. 7
- [51] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022. 3
- [52] Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal chain-of-thought reasoning in language models. *arXiv preprint arXiv:2302.00923*, 2023. 3, 7
- [53] Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibe Yang. Ddcot: Duty-distinct chain-of-thought prompting for multimodal reasoning in language models. *Advances in Neural Information Processing Systems*, 36:5168–5191, 2023. 2, 4, 7
- [54] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 3
- [55] Yuchen Zhuang, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. Toolqa: A dataset for llm question answering with external tools. *Advances in Neural Information Processing Systems*, 36, 2024. 3

Cantor[✂]: Inspiring Multimodal Chain-of-Thought of MLLM

Supplementary Material

A. Prompts used in Cantor

Cantor consists of two stages: Decision-Generation and Execution. For different models, Cantor (Gemini) and Cantor (GPT-3.5) use the same prompt during the Decision-Generation stage and the Execute-Synthesis stage. During the Execute-Synthesis stage, we use the generated sub-tasks as the prompts. In terms of prompts for different datasets, only the in-context learning examples are different, while the other prompts are used the same.

A.1. Decision-Generation

In the Decision-Generation stage, the prompt template we used includes task instructions and in-context learning examples. The task instructions, as shown in Fig 4, actually include the guidance for the task, the functional definition of the expert module, and the format requirements for the answers. We design different in-context learning examples based on the characteristics of ScienceQA and MathVista. For ScienceQA, we use in-context learning examples as shown in Figs 5, 6, and 7. For MathVista, we use in-context learning examples as shown in Figs 8 and 9.

A.2. Execute-Modularization

In the Execute-Modularization stage, Cantor executes the sub-tasks assigned by the Decision-Generation by calling the expert module played by MLLM. Therefore, at this stage, we do not use manually designed prompts and directly input the assigned sub-tasks into MLLM as prompts, in the specific format of [Expert Module: a corresponding sub-task]. For example: [ChartSense Expert: Extract the values of all the bars from the chart.] and [VisionIQ Analyst: What is the total number of people in the image?].

A.3. Execute-Synthesis

We use the prompt template shown in Fig 10 to generate answers during the Execute-Synthesis stage. This includes prompts for generating answers and formatting requirements for answers. This includes three key points: **1.** Play the role of an Answer Generator who is knowledgeable and adept at integrating information. **2.** Think carefully before answering. **3.** Maintain rationality and criticality when dealing with supplementary information.

B. Additional Analysis of Cantor

B.1. Case Presentation

In Figs 11, 12, 13, we show some specific cases of Cantor. It can be seen that Cantor has good decision-generating and practical problem-solving abilities.

B.2. Ablation Study with Modules

In this section, we further analyze the ablation study of the expert module. In the ablation experiment in the main text, we find that enabling a module has a greater impact on model performance than disabling a module. We speculate that this is because when MLLM acts as various experts, it possesses a certain degree of universal higher-level information capture capability. As shown in Fig 14, after disabling ChartSense Expert, Cantor will adaptively adjust the decision and instead ask VisionIQ Analyst questions to obtain information about the chart. And VisionIQ Analyst also correctly answers this sub-task and facilitates the final inference to obtain the correct answer. This case illustrates that thanks to the versatility of MLLM, when playing various experts using MLLM, even if one expert module is disabled, the remaining expert modules can to some extent compensate for the lack of that module.

However, disabling a certain expert module still affects the integrity of Cantor. As shown in Fig 15, for chart information extraction, compared to ChartSense Expert, when using VisionIQ Analyst, it only extracts data for three bars and ignores the other bars with values of 0. This indicates that although different expert modules are to some extent universal, they are not omnipotent. Specific expert modules still focus on specific abilities and are lacking in other professional abilities. This also demonstrates the importance and irreplaceability of the four expert modules we propose. At the same time, we believe that thanks to the excellent scalability of Cantor, introducing more expert modules with different functions will further improve its performance.

B.3. Impact of Visual Information Levels

In this section, we demonstrate the impact of different levels of visual information on Gemini’s decision generation. As shown in Fig 16, when asking which country is highlighted, the model cannot answer the question both in the absence of visual information and with only a rough caption provided. This is because the model cannot acquire effective visual information solely from the question or a rough caption. When entering a detailed caption, even if it contains a lot of content, it is irrelevant information and lacks

key information about what the highlighted country is. The model still cannot answer the question. Only by inputting images can the model obtain sufficient visual information for problem-solving.

As shown in Fig 17, another case is shown. When the detailed caption contains key information to answer the question, the model can also provide the correct answer. However, it should be noted that in practical applications, we cannot control whether the captions include key information for solving problems. On the contrary, the image must contain clues to the problem-solving. Therefore, inputting images is the best way to obtain visual information during decision generation.

Decision Generation Prompt

You are an advanced question-answering agent equipped with four specialized modules to aid in analyzing and responding to queries about images:

1. **TextIntel Extractor:** This module extracts and converts text within images into editable text format. It's particularly useful for images containing a mix of text and graphical elements. When this module is required, specify your request as: "TextIntel Extractor: <specific task or information to extract>."

2. **ObjectQuant Locator:** This module identifies and locates objects within an image. It's adept at counting objects and determining their spatial arrangement. When you need this module, frame your request as: "ObjectQuant Locator: <object1, object2, ..., objectN>," listing the objects you believe need detection for further analysis.

3. **VisionIQ Analyst:** This module processes and interprets visual data, enabling you to ask any queries related to the image's content. When information from this module is needed, phrase your request as: "VisionIQ Analyst: <your question about the image>."

4. **ChartSense Expert:** This module specializes in analyzing and interpreting information from charts and graphs. It can extract data points, understand trends, and identify key components such as titles, axes, labels, and legends within a chart. When you require insights from a chart or graph, specify your request as: "ChartSense Expert: <specific aspect of the chart you're interested in or question you have about the chart>."

When faced with a question about an image, which will be accompanied by a hint that might not cover all its details, your task is to:

If the question can be answered directly based on the information provided without the need for detailed input from the modules, specify this explicitly. Do not disclose the answer itself.

Otherwise:

- Provide a rationale for your approach to answering the question, explaining how you will use the information from the image and the modules to form a comprehensive answer.
- Assign specific tasks to each module as needed, based on their capabilities, to gather additional information essential for answering the question accurately.

Your response should be structured as follows:

Answer:

["This question does not require any modules and can be answered directly based on the information provided."] or [Rationale: Your explanation of how you plan to approach the question, including any initial insights based on the question and image information provided. Explain how the modules' input will complement this information.]

Modules' tasks (if applicable):

1. TextIntel Extractor: [Specify the text or information to be extracted from the image, if necessary.]
2. ObjectQuant Locator: [List the objects to be identified or counted in the image, if required.]
3. VisionIQ Analyst: [Pose any specific questions you have about the image that require deeper visual analysis, if applicable.]
4. ChartSense Expert: [Extract chart data or specify any questions about the chart, if required.]

Ensure your response adheres to this format to systematically address the question using the available modules or direct analysis as appropriate.

Here are some examples:

...

(See Figure 5, 6, 7, 8, 9)

...

Please refer to the prompts and examples above to help me solve the following problem: <Question>.

Figure 4. The prompt of the Decision-Generation stage.

In-context Learning Examples (ScienceQA)

Here are some examples:

"Question1": "Which solution has a higher concentration of blue particles?",
"Choices": [
"Solution B",
"neither; their concentrations are the same",
"Solution A"
],

Answer:

1. Concentration in a solution refers to the amount of a substance (solute) present in a specified amount of another substance (solvent).
2. To know the solvent volume, we need TextIntel Extractor to extract information about the volume from the image.
3. To understand the number of blue particles in solution A and solution B, we need a ObjectQuant Locator to detect them.

Modules' tasks:

1. TextIntel Extractor: Extract keywords related to solution volume in Solution A and Solution B.
2. ObjectQuant Locator: the number of blue particles in Solution A and Solution B.

"Question2": "Compare the average kinetic energies of the particles in each sample. Which sample has the higher temperature?",

"Choices": [
"neither; the samples have the same temperature",
"sample B",
"sample A"
],

Answer:

1. The temperature of particles in a substance is directly proportional to its average kinetic energy. The formula for kinetic energy of an particle is $\frac{1}{2}mv^2$, where m represents the mass of the object and v represents its velocity.
2. So TextIntel Extractor is needed to retrieve key information about mass and speed in the image.

Modules' tasks:

1. TextIntel Extractor: Extract keywords related to mass and speed in sample A and sample B.

Figure 5. In-context Learning Examples on ScienceQA.

In-context Learning Examples (ScienceQA)

"Question3": "Think about the magnetic force between the magnets in each pair. Which of the following statements is true?",

"Choices": [

"The magnetic force is stronger in Pair 2.",

"The magnetic force is stronger in Pair 1.",

"The strength of the magnetic force is the same in both pairs."

],

Answer:

1. The magnitude of magnetic force is independent of the direction of the magnetic pole, and is not related to whether the magnetic pole is attracted or repelled.

2. Different pairs of magnets will not affect each other.

3. The magnitude of magnetic force is inversely proportional to the distance between two magnets.

4. To know the distance between magnets, TextIntel Extractor is needed to extract distance information.

Modules' tasks:

1. TextIntel Extractor: Extract the distance values between two pairs of magnets in pair1 and pair2 separately.

"Question4": "What is the expected ratio of offspring with mutated antennae to offspring with normal antennae? Choose the most likely ratio.",

"Choices": [

"1:3",

"0:4",

"3:1",

"2:2",

"4:0"

]

Answer:

1. In a group of fruit flies, some individuals have mutated antennae and others have normal antennae. In this group, the gene for the antenna type trait has two alleles.

2. The allele 'A' is for mutated antennae, the allele 'a' is for normal antennae.

3. To know the specific gene composition, we need TextIntel Extractor to extract information from Punnett square.

Modules' tasks:

1. TextIntel Extractor: Extract the genotypes from the Punnett square.

Figure 6. In-context Learning Examples on ScienceQA.

In-context Learning Examples (ScienceQA)

"Question5": "Which month is the wettest on average in Christchurch?",

"Choices": [

"August",

"April",

"May"

]

Answer:

1. This question provides a chart of Christchurch precipitation. We need to compare the monthly precipitation in Christchurch to determine which month is the wettest.
2. In order to determine which month has the highest precipitation, we need ChartSense Expert to extract the precipitation for each month

Modules' tasks:

1. ChartSense Expert: Extract the precipitation for each month from the chart.

"Question6": "What is the capital of New Jersey?",

"Choices": [

"Augusta",

"Montpelier",

"Newark",

"Trenton"

]

Answer:

1. This question does not require any modules and can be answered directly based on the information provided.

Figure 7. In-context Learning Examples on ScienceQA.

In-context Learning Examples (MathVista)

Here are some examples:

"Question1": Subtract all large rubber spheres. Subtract all big shiny cylinders. How many objects are left?

Answer:

1. First, we need to know the quantity of all objects. Secondly, it is important to know the number of large rubber spheres and the number of big shiny cylinders. Then perform the subtraction operation again.
2. To know the number of all objects, we need a ObjectQuant Locator.
3. To know the number of large rubber spheres, we need a ObjectQuant Locator.
4. To know the number of big shiny cylinders, we need a ObjectQuant Locator.

Modules' tasks:

1. ObjectQuant Locator: the number of all objects.
2. ObjectQuant Locator: the number of large rubber spheres.
3. ObjectQuant Locator: the number of big shiny cylinders.

"Question2": "What is the age gap between these two people in image?"

Answer:

1. To determine the age gap between two individuals, it is first necessary to know who they are and their date of birth.
2. Calculate the difference in their birth dates to obtain the age difference.

Modules' tasks:

1. VisionIQ Analyst: Who are the two people in the picture and what is their date of birth?

"Question3": "What time is shown? Answer by typing a time word, not a number. It is () after eight.",

"choices": [

"half",

"quarter",

"o'clock",

"quarter to",

"quarter past"

]

Answer:

1. To know the specific time, it is necessary to determine which number the hour and minute of the clock point to respectively.
2. Combine the numbers pointed by the hour hand and minute hand to obtain the final time and fill in the blanks.

Modules' tasks:

1. VisionIQ Analyst: the number the hour hand is pointing to.
2. VisionIQ Analyst: the number the minute hand is pointing to.

Figure 8. In-context Learning Examples on MathVista.

In-context Learning Examples (MathVista)

"Question4": "Is Light Seafoam less than Dark Salmon?",

"choices": [

"yes",

"no"

]

Answer:

1. We need to know the numerical values corresponding to Light Seafoam and Dark Salmon to determine their quantity relationship.
2. To obtain the numerical values corresponding to Light Seafoam and Dark Salmon, we need ChartSense Expert.

Modules' tasks:

1. ChartSense Expert: the value of Light Seafoam.
2. ChartSense Expert: the value of Dark Salmon.

"Question5": "How many algorithms have accuracies higher than 9?",

Answer:

1. First, we need to know the accuracy of each algorithm.
2. Then we can compare the accuracy of each algorithm with 9.
3. Finally, we can count the number of algorithms that have accuracy higher than 9.
4. To know the accuracy of each algorithm, we need ChartSense Expert to extract the information from the chart.

Modules' tasks:

1. ChartSense Expert: Extract the accuracy of each algorithm from the chart.

Figure 9. In-context learning examples on MathVista.

Execute-Synthesis Prompt

You are a knowledgeable and skilled information integration science expert. Please gradually think and answer the questions based on the given questions, options, and supplementary information. Please note that we not only need answers, but more importantly, we need rationales for obtaining answers. Please combine your knowledge and supplementary information to obtain reasoning and answers. Please prioritize using your knowledge to answer questions. If unable to answer, maintain critical thinking and select effective information to assist you in selecting the most correct option as the answer. Furthermore, please do not rely solely on supplementary information, as the provided supplementary information may not always be effective. Please do not answer with uncertainty, try your best to give an answer.

The expected response format is as follows: **Rationale:**<rationale> **Answer:**<answer>.

Please answer the following case: <Question> , <Supplementary information>.

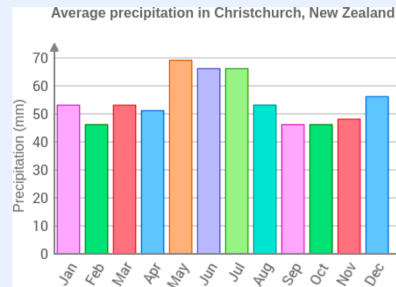
Figure 10. The prompt of the Execute-Synthesis stage.

The Specific Case of Cantor

Which month is the wettest on average in Christchurch?

Choices:

- (A) August
- (B) April
- (C) May



Cantor

Supplementary Information:

1. This question provides a chart of Christchurch precipitation. We need to compare the monthly precipitation in Christchurch to determine which month is the wettest.
2. In order to determine which month has the highest precipitation, we need ChartSense Expert to extract the precipitation for each month.

Modules' tasks:

ChartSense Expert: Extract the precipitation for each month from the chart.

Answer: The precipitation for each month is:\n- January: 52 mm\n- February: 46 mm\n- March: 51 mm\n- April: 50 mm\n- May: 68 mm\n- June: 67 mm\n- July: 67 mm\n- August: 52 mm\n- September: 47 mm\n- October: 48 mm\n- November: 49 mm\n- December: 54 mm.

Answer Generator:

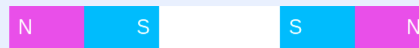
Rationale: The month with the highest precipitation is May, with 68 mm of rain.

Answer: May.

Will these magnets attract or repel each other?

Choices:

- (A) repel
- (B) attract



Cantor

Supplementary Information:

1. Like poles of a magnet repel each other. Unlike poles of a magnet attract each other.
2. To know the poles of the magnets, TextIntel Extractor is needed to extract information about the poles.

Modules' tasks:

TextIntel Extractor: Extract the poles of the two magnets.

Answer: The poles of the magnets are: - Magnet 1: North (N) and South (S) - Magnet 2: South (S) and North (N).

Answer Generator:

Rationale: The south pole of magnet 1 is facing the south pole of magnet 2. Like poles repel each other.

Answer: repel.

Figure 11. The Specific Case of Cantor.

The Specific Case of Cantor

What is the probability that a cow produced by this cross will have solid coloring?

Choices:

- (A) 4/4
- (B) 0/4
- (C) 1/4
- (D) 2/4
- (E) 3/4

| | | |
|---|----|----|
| | a | a |
| a | aa | aa |
| a | aa | aa |



Cantor

Supplementary Information:

1. The probability of getting a cow with solid coloring is determined by the number of offspring with the genotype AA or Aa.
2. To know the specific gene composition, we need TextIntel Extractor to extract information from Punnett square.

Modules' tasks:

TextIntel Extractor: Extract the genotypes from the Punnett square.

Answer: aa, aa, aa, aa

Answer Generator:

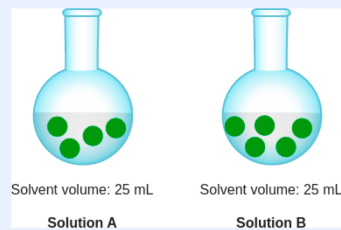
Rationale: Since all of the offspring have the aa genotype, this means that all of them will have white spots.

Answer: 0/4.

Which solution has a higher concentration of green particles?

Choices:

- (A) Solution B
- (B) neither; their concentrations are the same
- (C) Solution A



Cantor

Supplementary Information:

1. Concentration in a solution refers to the amount of a substance (solute) present in a specified amount of another substance (solvent).
2. To know the solvent volume, we need TextIntel Extractor to extract information about the volume from the image.
3. To understand the number of green particles in solution A and solution B, we need ObjectQuant Locator to detect them.

Modules' tasks:

TextIntel Extractor: Extract keywords related to solution volume in Solution A and Solution B.

Answer: Solution A: 25 mL Solution B: 25 mL

ObjectQuant Locator: the number of green particles in Solution A and Solution B.

Answer: Solution A has fewer green particles than Solution B.

Answer Generator:

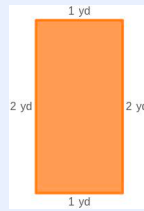
Rationale: Solution B has more green particles than Solution A.

Answer: Solution B.

Figure 12. The Specific Case of Cantor.

The Specific Case of Cantor

What is the perimeter of the rectangle?



Cantor

Supplementary Information:

To determine the perimeter of the rectangle, we need to know the length of all four sides and then add them together. Since the opposite sides of a rectangle are equal, we can calculate the perimeter using the formula: $\text{Perimeter} = 2(\text{length} + \text{width})$.

Modules' tasks:

VisionIQ Analyst: Determine the length and width of the rectangle.

Answer: In this case, VisionIQ Analyst would determine that the length of the rectangle is 2 yards and the width is 1 yard.

Answer Generator:

Rationale: To determine the perimeter of the rectangle, we need to know the length of all four sides and then add them together. Since the opposite sides of a rectangle are equal, we can calculate the perimeter using the formula: $\text{Perimeter} = 2(\text{length} + \text{width}) = 2(2+1) = 6$.

Answer: 6 yd.

What is the blue kite in the lower right corner shaped like?

Choices:

- (A) ferre
- (B) cat
- (C) cloud
- (D) octopus



Cantor

Supplementary Information:

1. To determine the shape of the blue kite in the lower right corner, we need to identify the blue kite in the lower right corner and then determine its shape.
2. To identify the blue kite in the lower right corner, we need to know the location of the blue kite in the lower right corner.
3. To determine the shape of the blue kite in the lower right corner, we need to know the shape of the blue kite in the lower right corner.

Modules' tasks:

ObjectQuant Locator: Identify the location of the blue kite in the lower right corner.

Answer: The blue kite in the lower right corner is located in the sky above the people.

VisionIQ Analyst: Determine the shape of the blue kite in the lower right corner.

Answer: The blue kite in the lower right corner is an octopus.

Answer Generator:

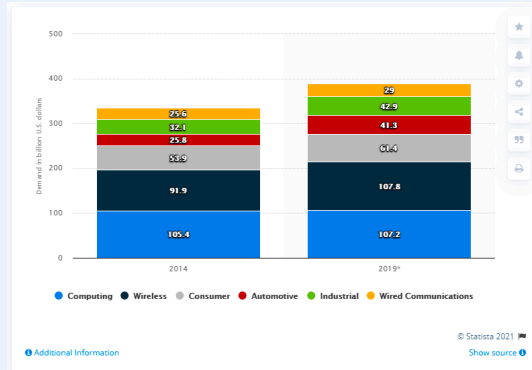
Rationale: The blue kite in the lower right corner is shaped like an octopus because it has eight tentacles.

Answer: (D)

Figure 13. The Specific Case of Cantor.

Ablation Study with Expert Module

What's the computing and wireless total for semiconductor demand in 2014?



Supplementary Information:

1. First, we need to know the computing demand and wireless demand in 2014. Then we can add them together to get the total.
2. To know the computing demand and wireless demand in 2014, we need to extract the information from the chart.

Modules' tasks:

ChartSense Expert: Extract the computing demand in 2014 from the chart.

Answer: The computing demand in 2014 was 105.4 billion dollars.

ChartSense Expert: Extract the wireless demand in 2014 from the chart.

Answer: 91.9 billion dollars.

Answer Generator:

Rationale: First, we need to know the computing demand and wireless demand in 2014. Then we can add them together to get the total. To know the computing demand and wireless demand in 2014, we need to extract the information from the chart. The computing demand in 2014 was 105.4 billion dollars. The wireless demand in 2014 was 91.9 billion dollars. So, the computing and wireless total for semiconductor demand in 2014 is $105.4 + 91.9 = 197.3$ billion dollars.

Answer: 197.3.



Cantor (without ChartSense Expert)

Supplementary Information:

1. First, we need to know the computing and wireless demand in 2014. Then we can add these two numbers together to get the total.
2. To know the computing demand in 2014, we need the VisionIQ Analyst to extract the information from the chart.
3. To know the wireless demand in 2014, we need the VisionIQ Analyst to extract the information from the chart.

Modules' tasks:

VisionIQ Analyst: What's the computing demand in 2014?

Answer: The computing demand for VisionIQ Analyst in 2014 was 105.4 billion US dollars.

VisionIQ Analyst: What's the wireless demand in 2014?

Answer: The wireless demand in 2014 was \$91.9 billion.

Answer Generator:

Rationale: First, we need to know the computing and wireless demand in 2014. Then we can add these two numbers together to get the total. The computing demand in 2014 was 105.4 billion US dollars. The wireless demand in 2014 was \$91.9 billion. So, the computing and wireless total for semiconductor demand in 2014 is $105.4 + 91.9 = 197.3$ billion US dollars.

Answer: 197.3.

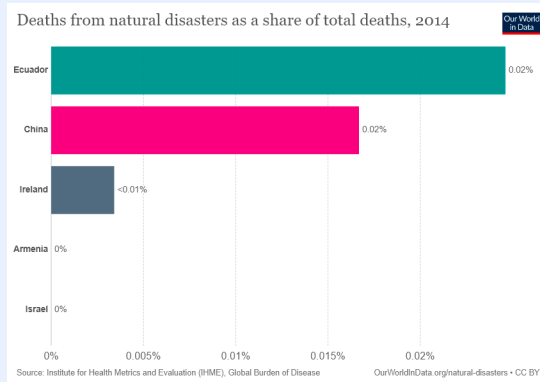
Figure 14. Case study on enabling and disabling Expert Module.

Ablation Study with Expert Module

Is the sum of two lowest bar is greater then the largest bar?

Choices:

- (A) Yes
- (B) No



Cantor

Supplementary Information:

1. First, we need to know the values of all the bars.
2. Then we can find the two lowest bars and add their values.
3. Finally, we can compare the sum of the two lowest bars with the value of the largest bar.
4. To know the values of all the bars, we need ChartSense Expert to extract the information from the chart.

Modules' tasks:

ChartSense Expert: Extract the values of all the bars from the chart.

Answer: Ecuador: 0.02% China: 0.02% Ireland: <0.01% Armenia: 0% Israel: 0%.

Answer Generator:

Rationale: First, we need to know the values of all the bars. We can use the ChartSense Expert to extract the information from the chart. The values of the bars are: Ecuador: 0.02%, China: 0.02%, Ireland: <0.01%, Armenia: 0%, Israel: 0%. Then we can find the two lowest bars and add their values. The two lowest bars are Armenia and Israel, and their values are 0% and 0%, respectively. So, the sum of the two lowest bars is $0\% + 0\% = 0\%$. Finally, we can compare the sum of the two lowest bars with the value of the largest bar. The largest bar is Ecuador, and its value is 0.02%. So, the sum of the two lowest bars is 0%, which is less than the value of the largest bar, 0.02%. Therefore, the answer is (B) No.

Answer: (B).



Cantor (without ChartSense Expert)

Supplementary Information:

1. To solve this problem, we need to know the values of all the bars.
2. Then we can calculate the sum of the two lowest bars and compare it with the largest bar.
3. Finally, we can determine whether the sum of the two lowest bars is greater than the largest bar.
4. To know the values of all the bars, we need VisionIQ Analyst to extract the information from the chart.

Modules' tasks:

VisionIQ Analyst: Extract the values of all the bars from the chart.

Answer: Ecuador: 0.02%\nChina: 0.02%\nIreland: <0.01%

Answer Generator:

Rationale: The sum of two lowest bar is $0.02\% \text{ (China)} + 0.01\% \text{ (Ireland)} = 0.03\%$. The largest bar is $0.02\% \text{ (Ecuador)}$. $0.03\% > 0.02\%$. So the answer is yes. Answer: (A).

Figure 15. Case study on enabling and disabling Expert Module.

Different levels of visual information

Which country is highlighted?

Choices:

- (A) Australia
- (B) Papua New Guinea
- (C) Fiji
- (D) New Zealand



No visual Information

Response:

1. Read the question: Which country is highlighted?
2. Look at the provided choices: ['Australia', 'Papua New Guinea', 'Fiji', 'New Zealand']
3. Since no visual information is provided, we cannot determine which country is highlighted.
4. Therefore, we cannot answer the question.

+ Rough Caption

Rough Caption: A blurry picture of a blue and white object.

Response:

1. The caption mentions a blurry picture of a blue and white object.
2. The choices are all countries.
3. We can't identify the object in the picture, so we can't use it to determine the country.
4. Therefore, we can't answer this question.

+ Detailed Caption

Detailed Caption: This is a map of Australia and New Zealand. Australia is the sixth-largest country in the world by land area, and New Zealand is the 75th-largest. Australia is located in the Southern Hemisphere and is bordered by the Indian Ocean to the west and the Pacific Ocean to the east. New Zealand is located in the South Pacific Ocean and is made up of two main islands, the North Island and the South Island.

Response:

1. Read the caption and find that the map shows Australia and New Zealand.
2. The question asks which country is highlighted.
3. The caption does not mention which country is highlighted.
4. The answer is not available in the context.

+ Image

Response:

The highlighted country is New Zealand.

The map shows the Australia and Oceania region. Australia is the largest country in the region and is located on the Australian continent. New Zealand is a country located in the South Pacific Ocean. It consists of two main islands, the North Island and the South Island.

Figure 16. The impact of different levels of visual information on model's responses.

Different levels of visual information

Which property do these three objects have in common?

- Choices:
 (A) rough
 (B) yellow
 (C) sticky



No visual Information

Response:

Read the question carefully and identify the three objects in question. The question does not specify what the three objects are, so we cannot answer this question from the provided context. The correct answer is: None of the above.

+ Rough Caption

Rough Caption: A series of photos showing different types of toys.

Response:

1. The first photo shows a toy car. The toy car is yellow.
 2. The second photo shows a toy ball. The toy ball is sticky.
 3. The third photo shows a toy brick. The toy brick is rough.
- The three objects have the property of being toys in common.
 The final answer is ['Toys'].

+ Detailed Caption

Detailed Caption: The image shows three items. A pair of yellow rain boots with blue edges and white soles, a yellow rubber duck with a blue beak and black eyes, and a yellow lemon. The items are all labeled with their names.

Response:

1. Identify the three objects in the image: - Rain boots - Rubber duck - Lemon
 2. Examine each object's properties: -Rain boots: Yellow, waterproof, flexible - Rubber duck: Yellow, smooth, buoyant -Lemon: Yellow, oval-shaped, acidic
 3. Find the common property: - All three objects share the color yellow.
- Therefore, the correct answer is yellow.

+ Image

Response:

Rationale: The three objects in the image are a pair of rain boots, a rubber duck, and a lemon. They all have the property of being yellow.
 Answer: yellow

Figure 17. The impact of different levels of visual information on model's responses.