

NPHardEval4V: A Dynamic Reasoning Benchmark of Multimodal Large Language Models

Lizhou Fan* **Wenyue Hua*** **Xiang Li*** **Kaijie Zhu**
 University of Michigan Rutgers University Shandong University Microsoft Research Asia

Mingyu Jin **Lingyao Li** **Haoyang Ling** **Jinkui Chi**
 Rutgers University University of Michigan University of Michigan University of Michigan

Jindong Wang **Xin Ma** **Yongfeng Zhang**
 Microsoft Research Asia Shandong University Rutgers University

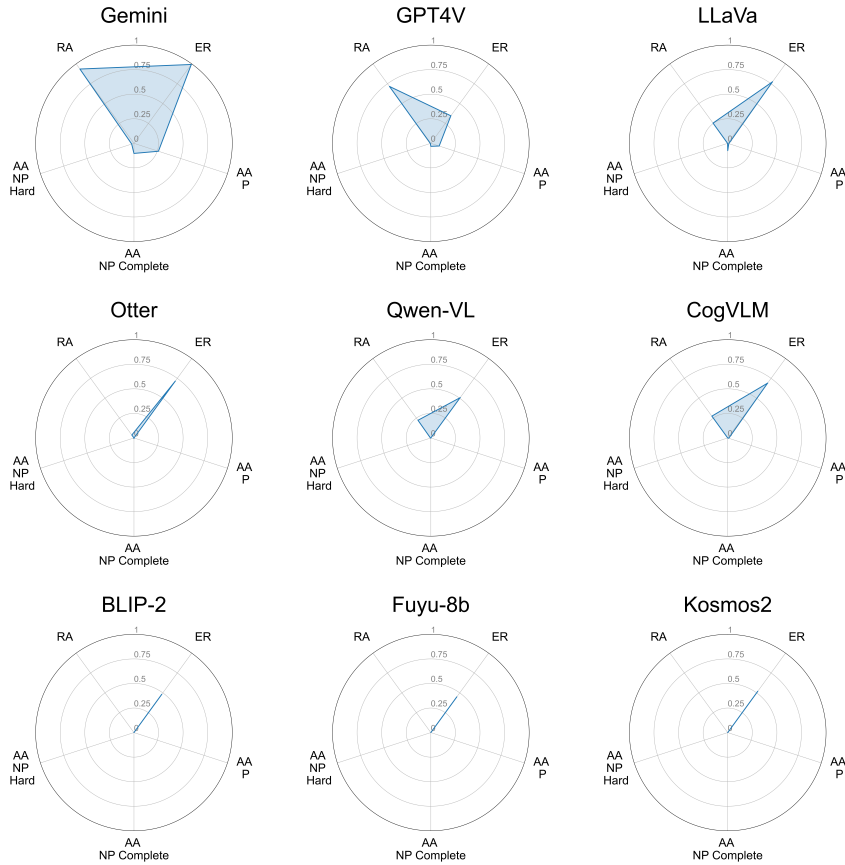


Figure 1: Multimodal Large Language Models’s performance on recognition (RA), Instruction-following (ER), and reasoning (AA) on polynomial time, NP-complete, and NP-hard problems.

Lizhou Fan, Wenyue Hua, and Xiang Li made equal contributions. **Correspondence to:** lizhouf@umich.edu, wenyue.hua@rutgers.edu, yongfeng.zhang@rutgers.edu. More author information are available in Appendix.

Abstract

Understanding the reasoning capabilities of Multimodal Large Language Models (MLLMs) is an important area of research. In this study, we introduce a dynamic benchmark, **NPHardEval4V**, aimed at addressing the existing gaps in evaluating the pure reasoning abilities of MLLMs. Our benchmark aims to provide a venue to disentangle the effect of various factors such as image recognition and instruction following, from the overall performance of the models, allowing us to focus solely on evaluating their reasoning abilities. It is built by converting textual description of questions from NPHardEval to image representations. Our findings reveal significant discrepancies in reasoning abilities across different models and highlight the relatively weak performance of MLLMs compared to LLMs in terms of reasoning. We also investigate the impact of different prompting styles, including visual, text, and combined visual and text prompts, on the reasoning abilities of MLLMs, demonstrating the different impacts of multimodal inputs in model performance. Unlike traditional benchmarks, which focus primarily on static evaluations, our benchmark will be updated monthly to prevent overfitting and ensure a more authentic and fine-grained evaluation of the models. We believe that this benchmark can aid in understanding and guide the further development of reasoning abilities in MLLMs. The benchmark dataset and code are available at <https://github.com/lizhouf/NPHardEval4V>.

1 Introduction

The evolution of Multimodal Large Language Models (MLLMs) marks a significant milestone in the pursuit of artificial general intelligence (AGI), following the advancement of Large Language Models (LLMs) [1, 2]. It introduces new capabilities for understanding and generating content that spans both text and visual inputs, contributing to enhanced multimedia interaction systems and sophisticated cross-modal decision-making tools [3, 4]. Reasoning is a critical ability for MLLMs as it is a fundamental aspect of problem-solving and task completion. The ability to reason enables MLLMs to understand complex relationships between different modalities, draw logical conclusions, and make informed decisions based on the information available.

The assessment of reasoning abilities in MLLMs is a crucial aspect of evaluating their overall performance and guiding the development of more advanced models. For instance, MLLMs can be tasked with reasoning about the minimum number of operations required to transform one string into another by analyzing the visual representation of the two strings and identifying the differences between them. Similarly, given a map, MLLMs can be tasked about reasoning on the optimal route between two locations on a map by analyzing the visual representation of the map and identifying the shortest path based on the available information. The ability to reason enables MLLMs to draw logical conclusions and make informed decisions based on the information available. In this paper, we introduce a benchmark that focuses solely on evaluating the reasoning abilities of various MLLMs. By assessing the extent of their reasoning capabilities, we aim to provide valuable insights into the strengths and limitations of current MLLMs and guide future research towards developing more advanced models with enhanced reasoning abilities. Our benchmark aims to disentangle the effect of various factors such as image recognition and instruction following, from the overall performance of the models, allowing us to focus solely on evaluating their reasoning abilities.

Numerous benchmarks have been developed to assess the capabilities of MLLMs in various domains, such as visual question answering [5, 6, 7, 8], Optical Character Recognition [9], robustness [10], hallucination [11], and holistic overall performance [12, 13, 14, 15]. However, despite the breadth of these evaluations, none of them specifically focus on assessing the pure reasoning abilities of MLLMs, leaving a significant gap in our understanding of their reasoning capabilities. Furthermore, many of these benchmarks are static in nature, making them prone to overfitting and limiting their effectiveness in providing a comprehensive view of MLLMs' abilities [16]. To address these limitations, there is a need for a dynamic benchmark that specifically targets the evaluation of MLLMs' reasoning abilities and updates regularly to prevent overfitting.

In response to these limitations, we introduce **NPHardEval4V**, a dynamic benchmark specifically designed to quantitatively and rigorously evaluate the reasoning ability of MLLMs across a diverse

set of tasks. This benchmark aims to provide a rigorous framework for assessing MLLMs’ reasoning, leveraging the computational complexity hierarchy to explore the depths of models’ reasoning abilities. The dynamic updating mechanism also makes the benchmark robust in a way that static benchmarks cannot. Moreover, the dynamic nature of NPHardEval4V — with datasets that evolve over time — mitigates the risk of overfitting, ensuring that assessments remain relevant and challenging. This is crucial for fostering models that are genuinely capable of learning and adapting, rather than merely optimizing for static benchmarks.

The proposed benchmark is developed based on the NPHardEval benchmark, as presented in [17]. The NPHardEval benchmark comprises nine types of algorithmic problems, categorized into three polynomial time problems, three NP-complete problems, and three NP-hard problems. Each algorithmic problem consists of 100 instances with varying difficulty levels. To enable a direct comparison between MLLMs and LLMs, we have retained the problems from NPHardEval and converted their textual descriptions into visual representations.

Our research questions are as follow:

1. **Reasoning Performance Evaluation of MLLMs:** This study seeks to understand the variation in reasoning abilities among different MLLMs and to identify the factors that contribute to any observed performance gaps. By analyzing and separating the influence of recognition and instruction following, we evaluate the core reasoning abilities of MLLMs vary with the nature of complexity and difficulty of problems. This investigation will provide insights into the relative strengths and weaknesses of MLLMs in reasoning tasks. We provide the model details of the MLLMs we evaluate in Table 2.
2. **Effects of Vision Input on MLLMs’ Performance:** We investigate how the inclusion of vision prompts affects MLLMs’ reasoning abilities compared to pure text inputs. Through ablation studies of three prompt types, including “Figures with Limited Instructional Text Setup” (figure+limited_text), “Text-only Setup” (full_text_only), and “Vision-rich-text Setup” (figure+full_text), this research will assess the impact of vision inputs on models’ reasoning performance and examine how the combination of vision and text inputs can affect MLLMs’ performance.

2 Related Work

2.1 Multimodal Large Language Models (MLLMs) and Their Reasoning Abilities

MLLMs [4] can process and interpret various multimodal data streams, including imagery and textual content [18]. As such, MLLMs are able to surpass singularly moded LLMs and unlock new avenues for performing real-world applications. As summarized by Yin et al. (2023) [3], MLLMs are more intelligent, user-centric, and holistic as compared to their LLM counterparts. MLLMs can mimic the way of how humans perceive the environment by assimilating multi-sensory inputs that can complement each other. They can also foster intuitive user interactions and communications. The scope of tasks that MLLMs can assist with is significantly broader compared to LLMs, reinforcing their versatility in applications [3] such as engineering [19] and healthcare [20].

Reasoning is one of the fundamental intelligent behaviors, essential for solving complex real-world tasks [21, 22]. However, even in text-centric settings, LLMs lack proper reasoning abilities, such as dealing with NP-hard (nondeterministic polynomial time) mathematical problems [17]. The quest to explore and improve reasoning capabilities of Strong AI particularly within MLLMs remains a persistent challenge and a pursuit of ongoing research [23, 24, 25].

In the domain of MLLMs, researchers have explored a range of techniques, such as instruction-tuning and prompt engineering, to enhance multimodal reasoning [4]. The practice of instruction tuning, vital for in-context learning (ICL), has emerged as one of the key techniques in enhancing the reasoning abilities of these models. For example, frozen LLM is an initial MLLM to showcase ICL ability [26]. Later, the Flamingo model demonstrated strong ICL capabilities with a more sophisticated LLM coupled with massive-scale image and text data for its pre-training phase [27]. Prompt engineering, on the other hand, has illustrated its effectiveness in improving the reasoning abilities of MLLMs [28, 29]. As summarized by Wang et al. (2024) [4], this approach entails a variety of strategic implementations, such as representation learning [28], exemplar generation [29], and model interactions [30].

2.2 Benchmarks of Multimodal Large Language Models (MLLMs)

As the reasoning ability of MLLMs continues to advance, benchmarks have become instrumental in evaluating their performance and identifying areas that require improvement. Wang et al. (2024) suggest that a robust multimodal reasoning benchmark must fulfill three key criteria: (1) the integration of multimodal information, (2) the categorization of reasoning, and (3) in-depth annotations of the reasoning steps [4].

Previous research has established numerous benchmarks to gauge the performance of MLLMs [31, 32, 13, 11, 33]. Among these, Fu et al. (2023) introduced an extensive MLLM evaluation benchmark to assess models’ perception and cognition skills across 14 distinct subtasks, including commonsense reasoning and code reasoning [31]. Later, Li et al. (2023) presented a benchmark called SEED-Bench, comprising 19,000 multiple-choice questions with precise human annotations. SEED-Bench evaluates MLLMs across 12 dimensions, capturing the ability to understand image and video modalities [13]. In a more recent study, Zhang et al. (2024) examined the self-consistency of MLLM responses in the presence of common corruptions. They created MMCBench, an extensive benchmark encompassing over 100 prominent LLMs. This benchmark assesses cross-modal interactions among text, images, and speech, incorporating four key generative tasks: text-to-image, image-to-text, text-to-speech, and speech-to-text [34].

Although existing benchmarks have evaluated the abilities of MLLMs across various dimensions, they do not provide a clear picture of the pure reasoning ability of MLLMs because factors such as recognition, knowledge amount, instruction following, and others are all combined in the presented performances in these benchmarks. Given that MLLMs’ general performance are intricately dependent on their recognition process and instruction-following ability, our study aims to factor out other factors in order to see the pure reasoning process and assess them of MLLMs. Moreover, these benchmarks lack dynamic updating mechanisms, which increases the risk of MLLMs becoming overfitted to these benchmarks and restricts their ability to accurately reflect the full range of MLLMs’ capabilities [16]. To address this research gap, our study proposes a dynamic assessment framework for evaluating the reasoning performance of MLLMs. In addition, we utilize the computational complexity hierarchy to rigorously assess the extent to which these MLLMs can achieve in reasoning tasks [17].

3 Benchmark Construction

As we mentioned above, NPHardEval4V is built upon NPHardEval [17], using the same set of problems while transformed the input type from textual to visual. This section outlines the transformation of the NPHardEval benchmark to suit the evaluation of MLLMs within NPHardEval4V, emphasizing the multimodal aspects of the tasks and the dynamic nature of the challenges.

3.1 NPHardEval Benchmark and Our Transformation

We explore and build upon the NPHardEval Benchmark [17], a framework that segments task complexity into three primary computational complexity classes: P (polynomial time), NP-complete (nondeterministic polynomial-time complete), and NP-hard. These classifications serve to delineate the intrinsic difficulty and the computational resources required for solving tasks within each class, showcasing an increasing order of complexity. NPHardEval4V also adopts the same three computational complexity levels (P, NP-Complete, and NP-Hard), transforming textual description of the problems into visual representations¹.

To provide a comprehensive overview of task complexity and difficulty within these classes, we introduce a hierarchical categorization illustrated in Table 1, where tasks are further divided into 10 progressive difficulty levels. This granular difficulty grading enables a nuanced assessment of model performance across a spectrum of computational challenges, thereby offering valuable insights into the capabilities and limitations of current models’ reasoning abilities.

To enable the direct comparison with the text-only prompt input in the NPHardEval benchmark, we transform the text-only questions and represent the data part with figures. We further illustrate these transformations, with the data of the problem presented both textually and visually.

¹Robustness experiment on the benchmark is performed in [17].

Complexity Class	Task
NP Hard (Most Complex)	Graph Coloring Problem Optimization Version (GCP) Traveling Salesman Problem Optimization Version (TSP) Meeting Scheduling Problem (MSP)
NP Complete	Knapsack Problem (KSP) Traveling Salesman Problem Decision Version (TSP-D) Graph Coloring Problem Decision Version (GCP-D)
P (Least Complex)	Shortest Path Problem (SPP) Edit Distance Problem (EDP) Sorted Array Search (SAS)

Table 1: Complexity classes and tasks

Graph Data Transformation The general construction of graph data problems involves providing MLLMs with both a textual description and a visual representation. For example, in the Graph Coloring Problem (see Figure 2), MLLMs are given a textual prompt alongside a figure depicting a graph, challenging them to reason across modalities. The figure is generated using Python based on the textual description of the graph coloring problem (GCP): it includes information about the nodes and the connections between them.

- 1 Graph coloring refers to the problem of coloring vertices of a graph in such a way that no two adjacent vertices have the same color.
- 2 There are 7 vertices 1 to 7 in a graph. You may use 3 colors with alphabets from A, B, C,... to color the graph. The graph is provided.

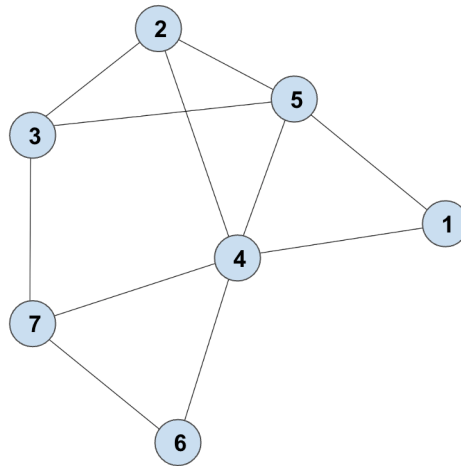


Figure 2: GCP Example

Figure 2 is transformed from the textual description of the below:

- 1 Vertex 1 is connect to 5, 4.
- 2 Vertex 2 is connect to 3, 4, 5.
- 3 Vertex 3 is connect to 5, 7.
- 4 Vertex 4 is connect to 1, 2, 6, 7.
- 5 Vertex 5 is connect to 1, 3, 4.
- 6 Vertex 6 is connect to 4, 7.
- 7 Vertex 7 is connect to 3, 4, 6.

The figures in algorithmic questions such as SPP, GCP_D, TSP, and TSP_D are transformed from text similarly: these figures include graphs containing nodes and edges; in the case of TSP and TSP_D, the weights on the edges are also included.

Linear Data Transformation Similarly, linear data problems like the Knapsack Problem (see Figure 3) are presented with both a prompt and a visual representation of items with associated weights and values. In order to provide a visual representation of the Knapsack problem, we use Python to create blocks with sizes corresponding to the weights of the items as described in the textual problem description. Specifically, each item is assigned a unique identifier (id) and a corresponding weight, which is used to determine the size of the block representing that item. The MLLMs are tasked with maximizing value within the constraints provided.

```
1 The Knapsack Problem (KSP) asks whether a subset of items, each with a
  given weight and value, can be chosen to fit into a knapsack of
  fixed capacity, maximizing the total value without exceeding the
  capacity. Determine if a subset of items can be selected to fit
  into a knapsack with a capacity of 40, maximizing weight without
  exceeding the capacity. Item weights are presented by the top
  number after 'W:'. Item IDs are shown as numbers below the weights
  after 'id:'.
```

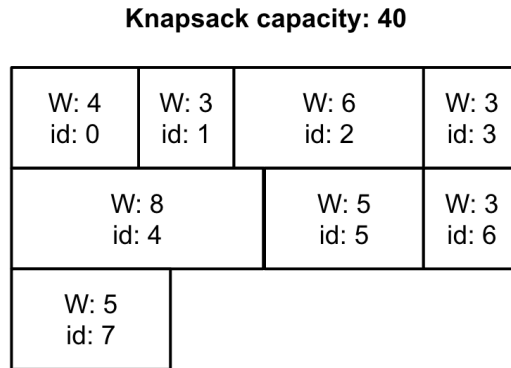


Figure 3: KSP Example

Figure 3 is transformed from the textual description of the below:

```
1 Item 0 has weight 4.
2 Item 1 has weight 3.
3 Item 2 has weight 6.
4 Item 3 has weight 3.
5 Item 4 has weight 8.
6 Item 5 has weight 5.
7 Item 6 has weight 3.
8 Item 7 has weight 5.
```

In the MSP problems, we generate figures that visually represent the availability of each person in a calendar format corresponding to the textual description of the problem. Specifically, we use Python to create a calendar view that displays the availability of each person based on the input data.

For algorithmic problems such as EDP and SAS, we do not employ a specialized visualization method. Instead, we simply display the relevant string of characters or string of numbers in the image.

To summarize, we employ various visualization functions to transform textual descriptions into visual representations for different algorithmic problems. In the prompt, we provide a combination of textual and visual information to the model, including an instructional prompt and a general introduction of the algorithm, along with an image that represents the specific problem to be solved. This approach is demonstrated in the previous examples for GCP and KSP. By providing both textual and visual information, we aim to evaluate the reasoning abilities of MLLMs in handling complex problems that require both visual and textual understanding.

4 Experimental Setting

To address the outlined research questions, our experimental setting systematically assesses the performance of various MLLMs in recognition and reasoning tasks and investigates the impact of vision and text inputs on MLLMs’ performance. The models we test on are presented in Figure 2. We detail the experimental setup corresponding to each research question, describing the models evaluated, the nature of the tasks, the methodology of the performance assessment, and the structure of the benchmark.

Model	Number of Parameters (B)	Access
GPT-4V [35]	Unknown	Close source
Gemini 1.0 Pro [36]	Unknown	Close source
CogVLM [37]	17	Open source
LLaVA-1.5-13B [38]	13	Open source
BLIP-2 FLAN-T5-XXL [39]	11	Open source
Fuyu-8B [40]	8	Open source
Otter [41]	8	Open source
Qwen-VL-7B [42]	7	Open source
Kosmos2 [43]	1.5	Open source

Table 2: MLLMs metadata

4.1 Recognition and Reasoning Performance Evaluation

In the initial experiment, the objective is to assess the performance of Multi-Modal Large Language Models (MLLMs) in terms of recognition and reasoning across a range of benchmark problems. To this end, zero-shot prompts are utilized to evaluate the inherent recognition and reasoning capabilities of each model. The performance of the models is then analyzed and compared across the benchmark problems to draw conclusions about their relative strengths and weaknesses in recognition and reasoning tasks.

Recognition Experiment In order to precisely assess the capacity for reasoning in Multi-Modal Large Language Models (MLLMs), it is essential to first comprehend their proficiency in image recognition, as accurate reasoning can only be established upon correct recognition. Therefore, an evaluation of recognition capabilities is initially conducted. For each question, both a visual representation of the description and a textual representation are provided. Subsequently, the MLLM is presented with both the visual and textual representations and is asked to determine whether they correspond to the same question. To mitigate the influence of randomness, the model is prompted five times, each time with a distinct random seed. The mean and variance of the model’s responses are then calculated, and it is concluded that the model demonstrates recognition capabilities if it correctly identifies the correspondence between the visual and textual representations more than half of the time, i.e., answering ‘yes’ three or more times. As an preprocess for analyzing reasoning ability, recognition experiment results and analysis are provided in Appendix B

Reasoning Experiment 1, the Default Setup (Vision with Instructional Text) This experiment evaluates the overall performance of MLLM on a benchmark consisting of various tasks. In the setting, a textual prompt is provided that includes a general introduction to the question and the format in which the answer should be given. Subsequently, an image is presented that pertains to the specific question being asked. The models are then tasked with processing both the textual and visual information in order to generate accurate responses. The performance of the MLLMs is evaluated based on their ability to correctly answer the questions, taking into account both their recognition and reasoning capabilities. The results of this experiment provide insights into the overall effectiveness of the models in handling multimodal tasks and their potential applications in real-world scenarios.

Below is an example for an SPP problem:

```
1 The Shortest Path Problem (SPP) involves finding the shortest path
   between two nodes in a weighted graph.
```

- 2 You need to find the shortest path between node 0 and node 3 in a graph. The graph's edges and their weights are given.
- 3 Please provide the shortest path from 0 to 3 and its total distance. Offer a concise step-by-step explanation of your reasoning process. Aim for brevity and clarity in your response.
- 4 Your output should be enclosed within `<root></root>` tags. Include your reasoning in `<reasoning></reasoning>` tags and the final path and total distance in `<final_answer></final_answer>` tags, like `<final_answer>{'Path': 'START->...->END', 'TotalDistance': 'INT_TOTAL_DISTANCE'}</final_answer>`.

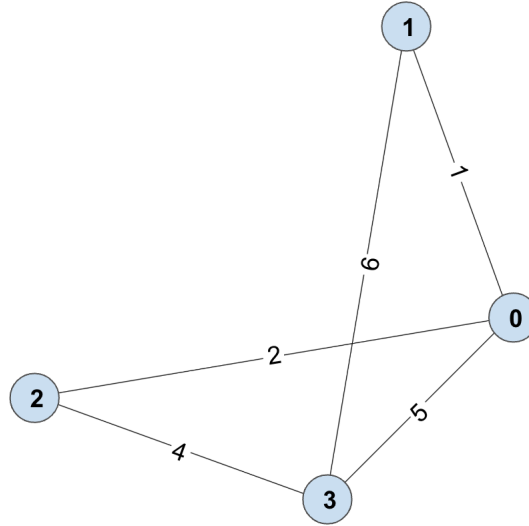


Figure 4: Image of an Short Path Problem example

In order to accurately evaluate the reasoning performance, it is necessary to isolate the specific aspect of reasoning ability from other factors that may influence the overall performance: we employ a filtering process that removes datapoints where the model fails to recognize the given input or fails to provide a parsable result, which indicates a failure in instruction-following. After filtering out these two factors, we can obtain a more accurate picture of the MLLMs' pure reasoning ability.

4.2 Ablation on Multimodal prompts

The goal of this experiment is to investigate the impact of prompt modality on the reasoning performance of MLLMs. Specifically, we aim to assess whether the inclusion of visual or textual inputs can enhance the models' problem-solving abilities. To achieve this, we conduct an ablation study, varying the modality of the prompts provided to the MLLMs by pure textual prompt and textual + visual prompt. By comparing the models' performance across different prompt modalities, we aim to identify which modality leads to the best performance and evaluate the contribution of each modality to the overall performance.

Reasoning Experiment 2, the Text-only Setup (Instructional and Data Text) In order to assess whether visual representations can be helpful, we first need to see the general problem-solving performance using pure textual prompt. In this experiment, we give MLLMs pure textual description of the problem containing the description of a specific question.

In the SPP example, the pure textual prompt is shown as below:

- 1 The Shortest Path Problem (SPP) involves finding the shortest path between two nodes in a weighted graph.
- 2 You need to find the shortest path between node 0 and node 3 in a graph. The graph's edges and their weights are given.


```

3 Please provide the shortest path from 0 to 3 and its total distance.
  Offer a concise step-by-step explanation of your reasoning process
  . Aim for brevity and clarity in your response.
4 Your output should be enclosed within <root></root> tags. Include your
  reasoning in <reasoning></reasoning> tags and the final path and
  total distance in <final_answer></final_answer> tags, like <
  final_answer>{'Path': 'START->...->END', 'TotalDistance': '
  INT_TOTAL_DISTANCE'}</final_answer>.
5
6 In this graph:
7 The distance between node 0 and 1 is 1,
8 the distance between node 1 and 3 is 9,
9 the distance between node 0 and 3 is 5,
10 the distance between node 0 and 2 is 2,
11 the distance between node 2 and 3 is 4.

```

Reasoning Experiment 3, the Vision-rich-text Setup (Vision with Instructional and Data Text)

Visual aids can be helpful in solving complex problems for humans. For instance, when dealing with a GCP, sketching a diagram of the graph and attempting different coloring schemes can aid in determining whether a valid solution exists. In this experiment, we aim to investigate whether providing MLLM with both textual and visual representations of a problem can improve their performance in solving it. To do so, we present the MLLMs with a complete textual description of the problem, followed by an accompanying image that depicts the problem visually. Our hypothesis is that the image can serve as a supplementary source of information to the textual description, potentially enhancing the model’s understanding of the problem and enabling it to generate more accurate solutions. By incorporating both textual and visual inputs, we aim to evaluate the MLLMs’ ability to integrate and process multimodal information in complementary with each other, which is a crucial aspect of their overall performance.

4.3 Evaluation Metrics

Our evaluation methodology incorporates a suite of metrics to assess the reasoning ability of MLLMs. Building upon the metrics established by [17], we expand our analysis to include a novel metric specifically tailored to quantify and rule out noise in analyzing models’ reasoning abilities. This sequence of metrics will start with the Recognition Accuracy (RA) metric, evaluating the recognition accuracy of the problems in the visual prompts, i.e., if the MLLMs can recognize what are in the questions. Among those questions that a MLLM can accurately recognize the input, we then use the Instruction-following Effective Rate (ER) metric, calculating the average chances that a MLLM cannot yield an output compatible with the rule-based answer parser. Finally, the Aggregated Accuracy (AA) will calculate the weighted accuracy of MLLMs’ answers correctness among those recognizable and parsable scenarios, aggregated with RA and ER. Figure 5 shows the sequence and relationship of these metrics.

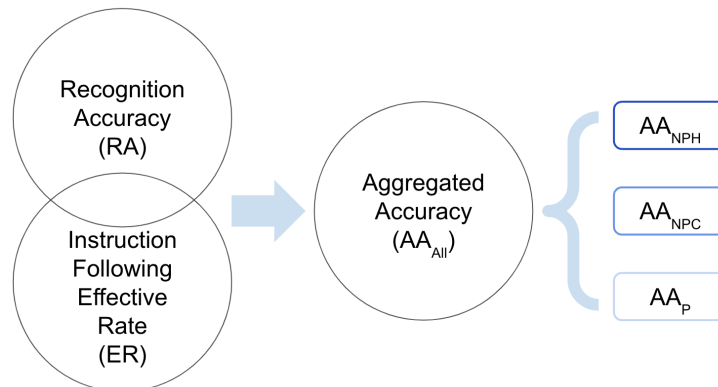


Figure 5: Metrics Sequence

Recognition Accuracy (RA) Recognition Accuracy evaluates the MLLMs’ capability to accurately interpret the visual information presented in the prompts. This metric is crucial as it forms the

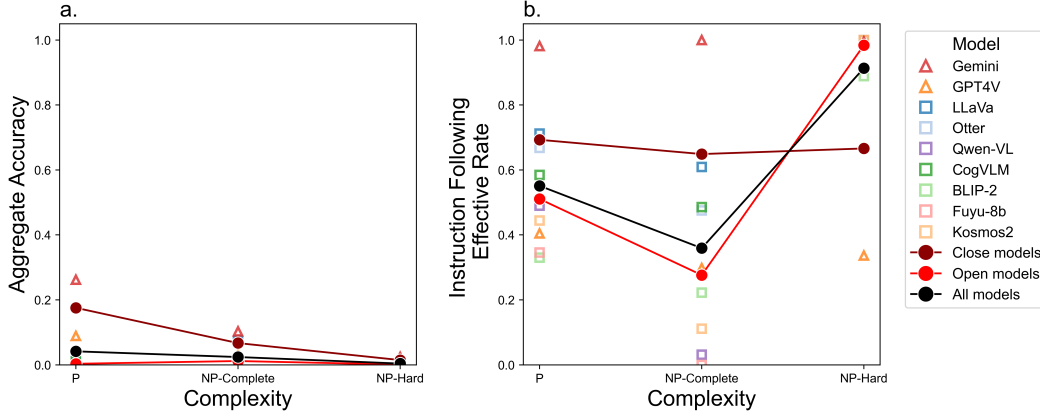


Figure 6: MLLM: a. Reasoning abilities performance excluding the effects of recognition and instruction following on figure+limited_text representation of questions. b. Instruction-following effective rate

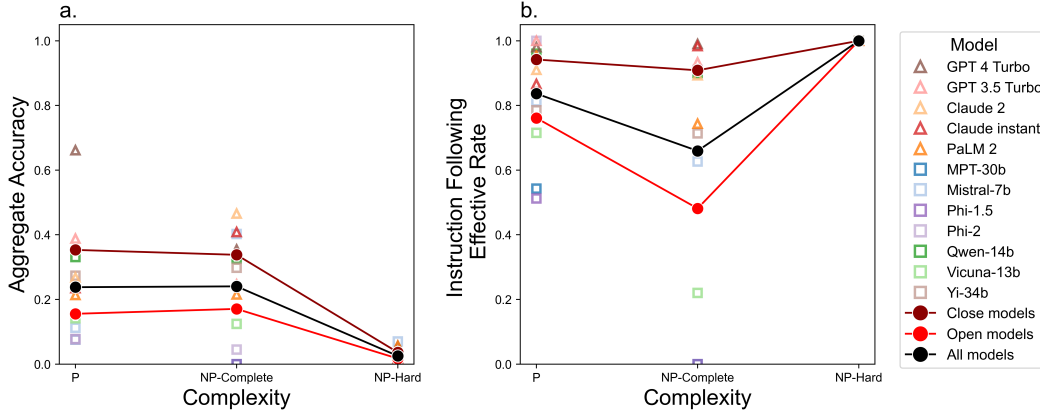


Figure 7: LLM: a. Reasoning abilities performance excluding the effect of instruction following on pure textual description of questions. b. Instruction-following effective rate

foundation for any subsequent reasoning or decision-making processes the model undertakes. It is quantified by the ratio of prompts that are correctly recognized by the model to the total number of prompts presented, and is defined as:

$$RA = \frac{\sum_{i=1}^N C_i}{N}$$

where C_i is a binary indicator of correct recognition (1 if the i^{th} prompt is correctly recognized, 0 otherwise), and N is the total number of prompts.

Instruction-following Effective Rate (ER) Instruction-following Effective Rate measures the average likelihood that an MLLM's response adhere to the expected output format, thus being compatible with a rule-based answer parser. This metric is crucial for gauging the models' reliability in producing solutions complying to standard parsing:

$$ER = \frac{\sum_{i=1}^N F_i}{N}$$

where F_i is parsing Instruction-following status for the i^{th} problem, where $F_i = 1$ if the parsing task successes and $F_i = 0$ otherwise, and N is the total number of problems.

Aggregated Accuracy (AA) Aggregated Accuracy takes into account the correct recognition of prompts and the successful parsing of responses to evaluate the correctness of MLLMs' answers.

This metric is a weighted measure that integrates the Recognition Accuracy and the inverse of the Failure Rate:

$$AA = \frac{\sum_{i=1}^N (w_i \times A'_i \times RA_i \times (ER_i))}{\sum_{i=1}^N w_i}$$

where w_i represents the difficulty weight for the i^{th} level, A'_i represents the accuracy in recognizable and parsable questions at i^{th} level, RA_i is the Recognition Accuracy for the i^{th} level, and ER_i is the Instruction-following Effectiveness Rate for the i^{th} level. **The Aggregated Accuracy metric will be used as the main evaluation metric for the overall performance and the complexity-specific submetrics, demoted as AA or AA_{All} , AA_{NPH} , AA_{NPC} , and AA_P .**

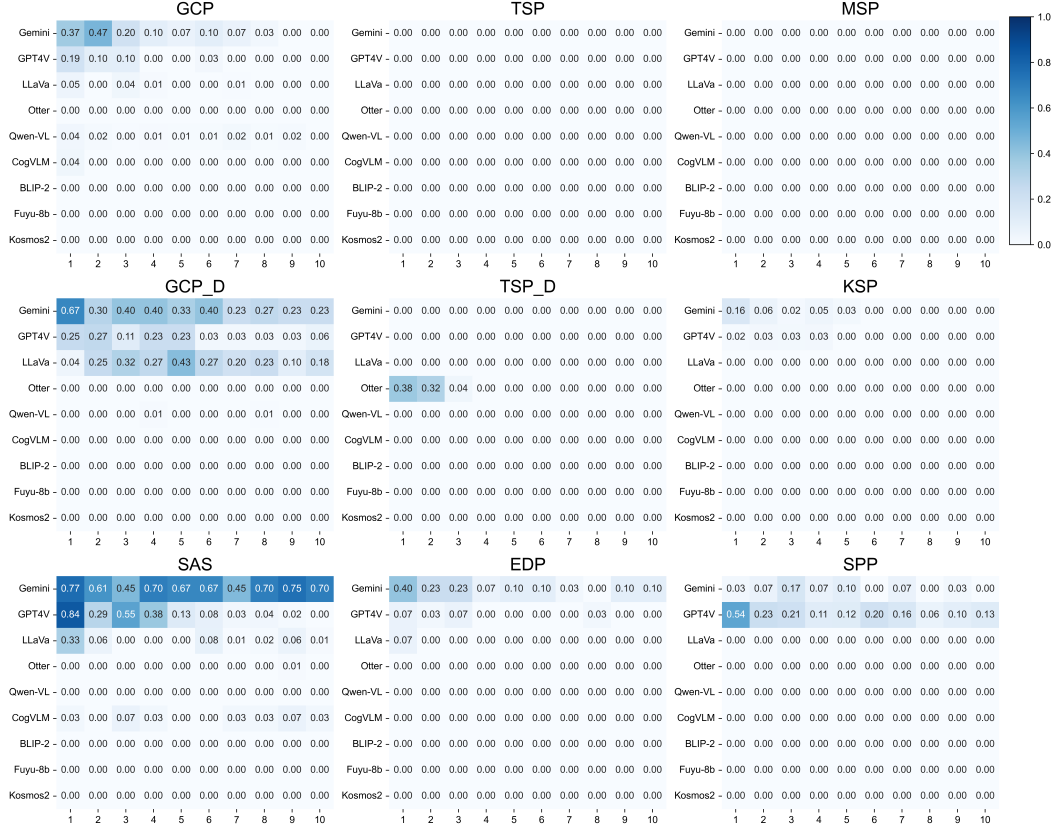


Figure 8: Reasoning abilities across models, complexity levels, and difficulty levels

5 Results

5.1 Reasoning Abilities of MLLMs

Close and Open Source Models As Figure 1 and Figure 7 indicate, the comparison between close source and open source MLLMs is quite stark, with close source models exhibiting superior performance in all tasks, irrespective of complexity class. Specifically, demonstrated in Figure 8, the Gemini model consistently outperforms its counterparts, including the other close source model GPT-4V, across the board. A detailed inspection of the heatmap reveals that Gemini maintains a higher aggregated accuracy and a more effective Instruction-following rate than GPT-4V, which suggests that the proprietary nature in the Gemini model may contribute to enhanced reasoning performance. This is especially evident in the NP-hard category, where Gemini shows remarkable resilience in reasoning accuracy compared to GPT-4V, which exhibits a notable decline.

Complexity Classes As Figure 8 indicates, the reasoning capabilities of MLLMs are inversely proportional to the complexity of the tasks. On simpler P problems, these models show commendable

performance. However, as the complexity escalates to NP-complete and further to NP-hard problems, a clear and expected downtrend in reasoning ability is observed. This trend is consistent across all MLLMs, signifying a universal challenge in tackling higher-order complexity with current multimodal reasoning architectures.

In addition to the overall trend, there are two notable findings. First, none of the nine models are capable of solving two of the NP-hard problems, TSP and MSP. Second, the Otter model performs unexpectedly well on the TSP-D problem – all other models, including the close source ones are unable to solve any of the TSP-D problems. This may due to the specific training data relevant or similar to the TSP-D problem.

Task Difficulties When focusing on individual reasoning tasks and considering models that can at least address the simplest questions, we notice a degradation in performance in correlation with increasing question difficulty. As Figure 8 shows, models like Gemini display a high success rate on easier questions within tasks like GCP, but this success rate diminishes as the difficulty level of the questions increases. This pattern underscores the models’ limitations and suggests that even the most capable MLLMs struggle to maintain their reasoning prowess as task complexity intensifies.

MLLMs vs. LLMs Figure 6 and Figure 7 provide a direct comparison between the current top-performing models in MLLMs and LLMs, including both closed-source and open-source models. In LLMs, we selected the models GPT-4-turbo, GPT-3.5-turbo, Claude-2, Claude-instant, PaLM-2, MPT-30b, Mistral-7b, Phi-1.5, Phi-2, Qwen-14b, and Yi-34b [17]. The results indicate that MLLMs lag behind LLMs in reasoning tasks. Specifically, the aggregated accuracy in P problems is approximately 0.4 in LLMs and 0.2 in MLLMs. Furthermore, the decrease in aggregated accuracy is more pronounced in MLLMs compared to LLMs. While LLMs maintain a relatively consistent weighted accuracy in NP-complete and P problems, MLLMs’ weighted accuracy decreases dramatically by half on average, from approximately 0.2 to 0.1. These findings highlight the need for further research and development to improve the reasoning abilities of MLLMs.

The results of our study indicate that the integration of multimodal data processing in MLLMs does not necessarily lead to improved reasoning capabilities. In fact, our findings reveal that the current development of MLLMs falls short in comparison to LLMs, with significantly weaker performance in reasoning tasks. This highlights the need for further research and development efforts aimed at enhancing the reasoning abilities of MLLMs.

5.2 Impact of Vision and Text Input

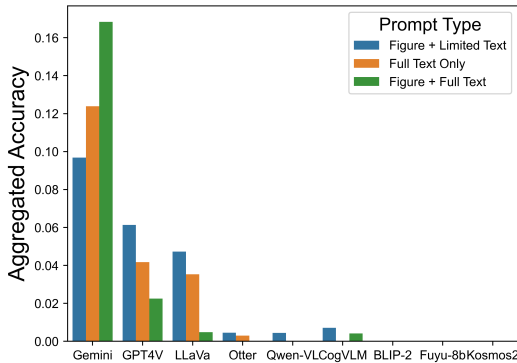


Figure 9: Reasoning abilities across prompt types

In evaluating the effect of different combinations of visual and textual inputs on the reasoning abilities of MLLMs, we observe varied responses across the models. Figure 9 indicates that a significant majority of the models, including both close source models like GPT-4V and open source models such as LLaVa, exhibit the highest levels of reasoning accuracy with the Default Setup, which pairs figures with limited instructional text (figure+limited_text). This suggests that their textual understanding is relatively weaker compared to their visual understanding, as the use of images to represent the question seems to be more effective than textual representation. Furthermore, when combining figure with full text, the models’ performance does not improve and instead decreases, indicating that the addition of excessive textual information may hinder their reasoning abilities. These findings

highlight the importance of considering the appropriate combination of visual and textual inputs in designing tasks for MLLMs, as well as the need for further research to enhance their textual understanding abilities.

In contrast, the Gemini model stands out distinctly by demonstrating superior reasoning performance with both the Text-only Setup (full_text_only) and the Vision-rich-text Setup (figure+full_text). This unique pattern indicates that Gemini may have a more advanced approach to processing and integrating comprehensive textual information, which is further enhanced by the presence of visual data. It's noteworthy that while the addition of full textual context to the visual prompts (Vision-rich-text Setup) improves Gemini's performance, it does not have the same positive effect on the other models tested, potentially due to the increased task complexity and vision recognition demands.

These findings reveal that the interplay between visual and textual inputs is not uniform across MLLMs and that certain models may be specially attuned to leverage text, visuals, or a combination of both to maximize their reasoning proficiency. It underscores the importance of prompt design in MLLM performance and calls for further investigation into how different models process multimodal information.

6 Conclusion and Discussion

In this paper, we have expanded upon the initial introduction of NPHardEval4V, a dynamic and comprehensive benchmark that scrutinizes the reasoning capabilities of Multimodal Large Language Models (MLLMs) against the backdrop of computational complexity. Our aim has been to dissect and understand the multifarious abilities of MLLMs, particularly focusing on their capacity for reasoning in response to visual-textual prompts.

The experimental results, complemented by the transformed data on the impact of vision and text input, suggest that the performance of MLLMs is contingent upon not only the complexity of the task but also the nature of input they are provided. We have observed that while most models, including close source models like GPT-4V and open source ones like LLaVa, exhibit optimal performance with the Default Setup (figure with limited instructional text), it is the Gemini model that stands out, showing a marked improvement in reasoning ability when provided with text-only and vision-rich-text (figure with full textual descriptions) prompts.

NPHardEval4V operates under the principle that benchmarks must be as dynamic as the models they assess. This is reflected in the benchmark's diverse suite of tasks, which compel MLLMs to demonstrate not just recognition but sophisticated reasoning and the ability to learn adaptively. This benchmark highlights the critical need for MLLMs to process and learn from both visual and textual data effectively, underscoring the importance of dynamic evaluation tools that can truly evaluate MLLMs' progression.

In conclusion, the insights garnered from NPHardEval4V shed light on the present competencies and constraints of MLLMs. While we have detailed the specific outcomes of our experiments, the broader trends emphasize the necessity of dynamic and stringent testing to deepen our comprehension and further the advancement of AI. As we continue to extend the possibilities of MLLMs, benchmarks like NPHardEval4V will be instrumental in steering the evolution of models that are not just powerful but also multifaceted, adaptable, and truly intelligent. This study acts as a clarion call for the AI community to recognize the complexities of multimodal reasoning and to strive for models that can seamlessly integrate diverse inputs to reason and learn more like humans do.

6.1 Limitations

Inherent Bias in Model Comparisons The benchmarking process inherently favors models like Gemini over others due to the varied architectural strengths, which may not be fully captured in a uniform evaluation setup. Despite our attempts to level the playing field, certain models are naturally more adept at specific types of reasoning tasks due to their difference in training data, which could skew the overall comparative analysis. For example, Gemini's superior performance in text-only setups may reflect an underlying bias in the benchmark towards linguistic reasoning over multimodal reasoning.

Limited Scope of Reasoning Tasks Our benchmark only scratches the surface of the vast landscape of reasoning. The tasks chosen, while diverse, cannot possibly encompass the full spectrum of reasoning abilities that MLLMs might possess. There is a risk that the benchmark may overrepresent certain reasoning styles while underrepresenting others, leading to an incomplete picture of models’ true reasoning capacities.

Prompt Dependence The study highlights the variability in model performance based on prompt type, suggesting a significant dependence on how questions are framed. This raises concerns about the models’ robustness and their ability to generalize across different input conditions. It is unclear whether the observed performance reflects true reasoning ability or an artifact of the models’ sensitivity to specific prompt structures.

Multimodal Integration Challenges The weaker performance of MLLMs in comparison to LLMs suggests that the integration of multimodal data may not always be beneficial for reasoning tasks. This observation implies limitations in how current models process and integrate multimodal information, potentially indicating a misalignment between model architecture and task requirements.

6.2 Research Outlook

Longitudinal Learning Studies Investigating models’ learning curves over extended periods could offer valuable insights into their potential for growth and adaptation. Long-term studies that track how models evolve with additional training data or through incremental learning can provide a deeper understanding of the underlying mechanisms that drive progress in MLLMs and identify the key factors that contribute to long-term success.

Expanding Reasoning Taxonomies Research should continue to expand the taxonomy of reasoning tasks within benchmarks to ensure a comprehensive evaluation of MLLMs. This includes not only increasing the variety of tasks but also ensuring that they are representative of the diverse ways in which reasoning is applied in real-world contexts. Developing benchmarks that can simulate more complex, real-life reasoning scenarios will be crucial for advancing the field.

Harmonizing Model Evolution with Benchmarks Balancing the dynamic nature of benchmarks with the pace of model development is essential. Researchers could explore introducing phased or tiered updates to benchmarks that allow for more strategic model improvements, rather than continuous, potentially disjointed, monthly updates. This approach could better align the benchmark evolution with the natural R&D cycles in AI development.

References

- [1] Lizhou Fan, Lingyao Li, Zihui Ma, Sanggyu Lee, Huizi Yu, and Libby Hemphill. A bibliometric review of large language models research from 2017 to 2023. *arXiv preprint arXiv:2304.02020*, 2023.
- [2] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.
- [3] 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.
- [4] Yiqi Wang, Wentao Chen, Xiaotian Han, Xudong Lin, Haiteng Zhao, Yongfei Liu, Bohan Zhai, Jianbo Yuan, Quanzeng You, and Hongxia Yang. Exploring the reasoning abilities of multimodal large language models (mllms): A comprehensive survey on emerging trends in multimodal reasoning. *arXiv preprint arXiv:2401.06805*, 2024.
- [5] 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.
- [6] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering.

- In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.
- [7] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *Proceedings of the IEEE/cvf conference on computer vision and pattern recognition*, pages 3195–3204, 2019.
 - [8] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019.
 - [9] Yuliang Liu, Zhang Li, Hongliang Li, Wenwen Yu, Mingxin Huang, Dezhi Peng, Mingyu Liu, Mingrui Chen, Chunyuan Li, Lianwen Jin, et al. On the hidden mystery of ocr in large multimodal models. *arXiv preprint arXiv:2305.07895*, 2023.
 - [10] Yunqing Zhao, Tianyu Pang, Chao Du, Xiao Yang, Chongxuan Li, Ngai-Man Man Cheung, and Min Lin. On evaluating adversarial robustness of large vision-language models. *Advances in Neural Information Processing Systems*, 36, 2024.
 - [11] Fuxiao Liu, Tianrui Guan, Zongxia Li, Lichang Chen, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. Hallusionbench: You see what you think? or you think what you see? an image-context reasoning benchmark challenging for gpt-4v (ision), llava-1.5, and other multi-modality models. *arXiv preprint arXiv:2310.14566*, 2023.
 - [12] 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*, 2023.
 - [13] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023.
 - [14] 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*, 2023.
 - [15] 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. *arXiv preprint arXiv:2311.16502*, 2023.
 - [16] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 2023.
 - [17] Lizhou Fan, Wenyue Hua, Lingyao Li, Haoyang Ling, and Yongfeng Zhang. Nphardeval: Dynamic benchmark on reasoning ability of large language models via complexity classes. *arXiv preprint arXiv:2312.14890*, 2023.
 - [18] Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and S Yu Philip. Multimodal large language models: A survey. In *2023 IEEE International Conference on Big Data (BigData)*, pages 2247–2256. IEEE, 2023.
 - [19] Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zichong Yang, Kuei-Da Liao, et al. A survey on multimodal large language models for autonomous driving. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 958–979, 2024.
 - [20] Ankit Pal and Malaikannan Sankarasubbu. Gemini goes to med school: Exploring the capabilities of multimodal large language models on medical challenge problems & hallucinations. *arXiv preprint arXiv:2402.07023*, 2024.
 - [21] Fei Yu, Hongbo Zhang, and Benyou Wang. Nature language reasoning, a survey. *arXiv preprint arXiv:2303.14725*, 2023.
 - [22] Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403*, 2022.
 - [23] Wei-Ge Chen, Irina Spiridonova, Jianwei Yang, Jianfeng Gao, and Chunyuan Li. Llava-interactive: An all-in-one demo for image chat, segmentation, generation and editing. *arXiv preprint arXiv:2311.00571*, 2023.

- [24] Zhaoyang Liu, Zeqiang Lai, Zhangwei Gao, Erfei Cui, Zhiheng Li, Xizhou Zhu, Lewei Lu, Qifeng Chen, Yu Qiao, Jifeng Dai, et al. Controllm: Augment language models with tools by searching on graphs. *arXiv preprint arXiv:2310.17796*, 2023.
- [25] Ran Gong, Qiuyuan Huang, Xiaojuan Ma, Hoi Vo, Zane Durante, Yusuke Noda, Zilong Zheng, Song-Chun Zhu, Demetri Terzopoulos, Li Fei-Fei, et al. Mindagent: Emergent gaming interaction. *arXiv preprint arXiv:2309.09971*, 2023.
- [26] Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems*, 34:200–212, 2021.
- [27] 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.
- [28] Jing Yu Koh, Ruslan Salakhutdinov, and Daniel Fried. Grounding language models to images for multimodal inputs and outputs. 2023.
- [29] Jiaxian Guo, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Boyang Li, Dacheng Tao, and Steven Hoi. From images to textual prompts: Zero-shot visual question answering with frozen large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10867–10877, 2023.
- [30] Zhenhailong Wang, Manling Li, Ruochen Xu, Luowei Zhou, Jie Lei, Xudong Lin, Shuohang Wang, Ziyi Yang, Chenguang Zhu, Derek Hoiem, et al. Language models with image descriptors are strong few-shot video-language learners. *Advances in Neural Information Processing Systems*, 35:8483–8497, 2022.
- [31] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.
- [32] Chaoyou Fu, Renrui Zhang, Haojia Lin, Zihan Wang, Timin Gao, Yongdong Luo, Yubo Huang, Zhengye Zhang, Longtian Qiu, Gaoxiang Ye, et al. A challenger to gpt-4v? early explorations of gemini in visual expertise. *arXiv preprint arXiv:2312.12436*, 2023.
- [33] Yunkang Cao, Xiaohao Xu, Chen Sun, Xiaonan Huang, and Weiming Shen. Towards generic anomaly detection and understanding: Large-scale visual-linguistic model (gpt-4v) takes the lead. *arXiv preprint arXiv:2311.02782*, 2023.
- [34] Jiawei Zhang, Tianyu Pang, Chao Du, Yi Ren, Bo Li, and Min Lin. Benchmarking large multimodal models against common corruptions. *arXiv preprint arXiv:2401.11943*, 2024.
- [35] OpenAI. Gpt-4v(ision) system card, 2023.
- [36] Google. Meet the first version of gemini— our most capable ai model., 2023.
- [37] Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023.
- [38] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.
- [39] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023.
- [40] Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşırklar. Introducing our multimodal models, 2023.
- [41] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Mimic-it: Multi-modal in-context instruction tuning. *arXiv preprint arXiv:2306.05425*, 2023.
- [42] 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.

[43] 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 preprint arXiv:2306.14824*, 2023.

A Author Information

Author Affiliations: L. Fan, L. Li, H. Ling, J. Chi: School of Information, University of Michigan, Ann Arbor, MI 48103, US; W. Hua, M. Jin, Y. Zhang: Department of Computer Science, Rutgers University, New Brunswick, NJ 08854, US; X. Li, X. Ma: School of Control Science and Engineering, Shandong University, Jinan PRC 250061; K. Zhu, J. Wang: Microsoft Research Asia, Beijing PRC 100080.

Author Emails: lizhouf@umich.edu, wenyue.hua@rutgers.edu, lixiang0814@mail.sdu.edu.cn, v-zhukaijie@microsoft.com, u9o2n2@u.northwestern.edu, lingyaol@umich.edu, hyfrankl@umich.edu, chijk@umich.edu, jindong.wang@microsoft.com, maxin@sdu.edu.cn, yongfeng.zhang@rutgers.edu

B Recognition Accuracy

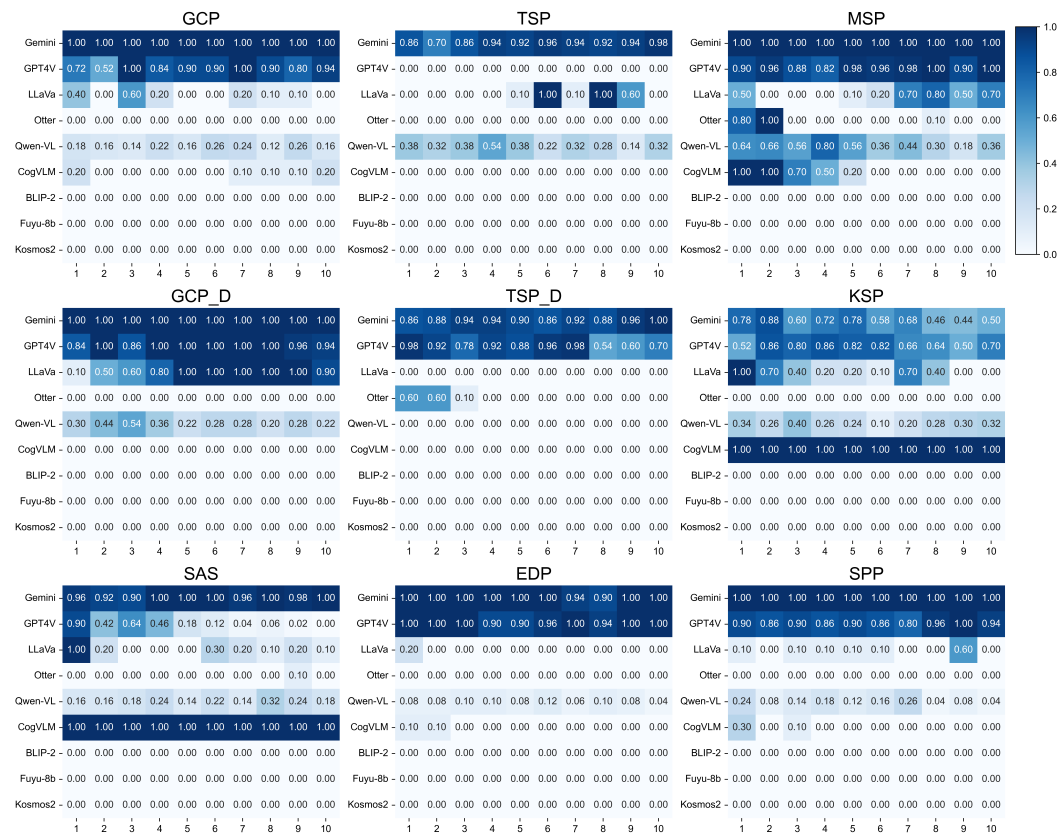


Figure 10: Prompt recognition across models, complexity levels, and difficulty levels. Recognition is an important preprocess that disentangles input processing’s impact on reasoning ability of MLLMs. In general, close source models outperforms open source counterparts in recognition. It is also important to notice irrelevance between complexity of tasks and recognition accuracy.

Several close source models, however, can obtain similar or even better recognition ability than the open source ones. For instance, CogVLM constantly reaches full recognition rate (1.0) on SAS and KSP problems. LLaVa, Otter, and QWen-VL also have different extents of outstanding recognition abilities on different tasks. In addition to further showing MLLMs’ variation in tasks regarding recognition ability, this result also reveals the individual strengths of different MLLMs, possibly due to different training data and techniques.