# CONSTRUCTURE: Benchmarking CONcept STRUCTUre REasoning for Multimodal Large Language Models

**Anonymous ACL submission**

## Abstract

Multimodal Large Language Models (MLLMs) have shown promising results in various tasks, but their ability to perceive the visual world with deep, hierarchical understanding similar to humans remains uncertain. To address this gap, we introduce CONSTRUCTURE, a novel concept-level benchmark to assess MLLMs' hierarchical concept understanding and reasoning abilities. Our goal is to evaluate MLLMs across four key aspects: 1) Understanding atomic concepts at different levels of abstraction; 2) Performing upward abstraction reasoning across concepts; 3) Achieving downward concretization reasoning across concepts; and 4) Conducting multi-hop reasoning between sibling or common ancestor concepts. Our findings indicate that even state-of-the-art multimodal models struggle with concept structure reasoning (e.g., GPT-4o averages a score of 62.1%). We summarize key findings of MLLMs in concept structure reasoning evaluation. Morever, we provide key insights from experiments using CoT prompting and fine-tuning to enhance their abilities.

## 1 Introduction

> *The basic level is the level in a taxonomy at which things are normally named, in the absence of reasons to the contrary. 'Dog' is a basic level category, 'boxer' a subordinate category, 'quadruped' a superordinate category.*
>
> — John R. Taylor

According to the prototype theory (Taylor, 2019) in cognitive science, humans perceive the visual world hierarchically, with basic, subordinate, and superordinate categories. People interpret the world differently based on these conceptual levels in diverse environments. As shown in Figure 1, humans can not only understand basic concepts,
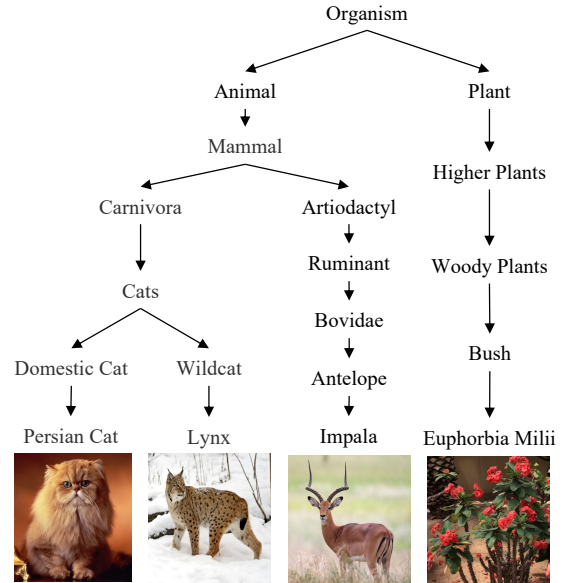


Figure 1: Demo of Concept Structure.

but also have a deep concept structure cognition in minds (Medin, 1989; Suresh et al., 2023). This raises the question: do multimodal AI systems, like Multimodal Large Language Models(*e.g.* GPT-4V (OpenAI, 2023)), exhibit similar concept structure cognition ability?

Recently, pretrained multimodal large language models (MLLMs) have transcended the confines of text-only modalities, gaining a deeper conceptual understanding of the world and demonstrating remarkable performance across a spectrum of downstream tasks. As a result, there is a growing importance and challenge in conducting comprehensive evaluations of these models to better understand their capabilities and pinpoint areas for enhancement. Inspired by the cognitive processes underlying human conceptual understanding, we posit that similar deep and structured visual conceptual cognition capabilities are pivotal for the profound comprehension of MLLMs. This compelling question demands exploration: Do MLLMs truly grasp

| Benchmark | Assessment of cognitive related abilities | Data Source | Answer Type | Evaluator | Size |
|---|---|---|---|---|---|
| LVLM-eHub (Xu et al., 2023) | Visual Reasoning, Visual Commonsense | Datasets | MC/OE | Metrics/LLMs/User | 332k |
| MME (Fu et al.) | Commonsense Reasoning, Numerical Calculating, Text Translation, Code Reasoning | Datasets | MC | Accuracy | 2,194 |
| MMBench (Liu et al., 2023) | Reasoning(*e.g.* Phsycial Relation Reasoning, Social Relation Reasoning) | Datasets/Handcraft/LLMs | MC | Accuracy | 2,974 |
| SEED-Bench (Li et al., 2023a) | Visual Reasoning, Spacial Relation | Handcraft/LLMs | MC | Accuracy | 19k |
| MM-Vet (Yu et al., 2023) | Spatial awareness, Knowledge, Math | Datasets/Handcraft | OE | LLMs | 218 |
| EgoThink (Cheng et al., 2023) | Scene Reasoning, Planning | Handcraft | OE | LLMs | 700 |
| **CONSTRUCTURE**(Ours) | **Concept Structure Reasoning** | **Datasets/Handcraft** | **MC** | **Accuracy** | **2,064** |

Table 1: Comparison of recent comprehensive evaluation benchmarks of MLLMs and our proposed benchmark $M^2C^2$-Bench.

and internalize concepts, or do they merely learn the superficial concept alignment through pretraining?

However, addressing this question requires a deeper investigation into the underlying mechanisms and limitations of MLLMs. As shown in Table 1, the deep visual concept structure cognition remains largely unexplored in existing benchmarks. On one hand, these benchmarks seldom consider visual cognitive capabilities. On the other hand, existing benchmarks for visual cognition typically focus solely on assessing visual reasoning ability related to the whole image content, overlooking the evaluation of deep and structural visual concept cognition. For example, SEED-bench **??** evaluates spatial relationships or visual reasoning tasks based on the content of images, while MME **??** assesses common-sense reasoning, numerical computation, code inference, and text translation. LVLM-eHub **??** evaluates visual reasoning and visual common sense. However, all these benchmarks require answering questions based on the entire content of an image. Therefore, there is a pressing need to construct an benchmark specifically designed to assess the deep visual concept structure cognition in MLLMs.

In this work, we propose the novel CONSTRUCTURE benchmark focusing on deep visual concept structure cognition. To fully uncover the deep visual cognitive ability of MLLMs across the overall concept structure beyond the Figure 1, we consider the following four key capabilities, ranging from *atomic concenpt understanding*, *concept abstraction reasoning*, *concept concretization reasoning*, to *common ancestor reasoning*. The first capability is to evaluate the atomic concept understanding at different levels. The last three capabilities are to evaluate reasoning ability on concept structure, including upward abstraction reasoning from child concepts to parent concepts, downward concretization reasoning from parent concepts to child concepts and multi-hop reasoning between sibling concepts or common ancestor concepts, respectively.

Based on the our proposed CONSTRUCTURE benchmark, we conduct comprehensive experiments to evaluate concept structure cognition capabilities of fourteen popluar MLLMs (including 6 api-based MLLMs and 8 open-sourced MLLMs). We conclude the main findings as follows:

1) Current MLLMs possess a certain level of conceptual understanding, but their performance in concept structure reasoning is poor. The best model, GPT-4o, only achieved a score of 0.621, indicating significant room for improvement.

2)In concept structure reasoning tasks, MLLMs perform the worst in common ancestor reasoning. The main reasons for errors are inconsistencies in the reasoning process and constraint violations. This demonstrates that adhering to multiple constraints and maintaining consistency in the reasoning process are key challenges to enhancing MLLMs' concept structure cognitive abilities.

3) MLLMs still need to improve their ability to reason about hierarchical relationships in concept structures.

4) MLLMs' performance deteriorates as the concept hierarchy deepens and granularity increases. Therefore, improving fine-grained concept recognition and the ability to reason about related fine-grained concept structures is crucial for enhancement.

Our evaluation results reveal the limitations of MLLMs in concept structure cognition, providing a comprehensive and clear analysis that directs further improvements for MLLMs. Furthermore, we improved MLLMs' concept structure reasoning abilities through few-shot CoT prompting and fine-tuning methods, and elucidated three key insights in the discussion.

## 2 Related Work

### 2.1 Multimodal Large Language Models

Building upon the monumental achievements of large language models (LLMs) (Du et al., 2022; OpenAI, 2023; Zhu et al., 2023), recent advance-

ments in multimodal LLMs (MLLMs) have firmly established LLMs as their foundational backbone. Initially, MLLMs commence pre-training using large-scale image-text paired datasets (Yang et al., 2023; Li et al., 2023b) or by integrating random visual and textual data (Ye et al., 2023; Bai et al., 2023). This foundational phase is further enriched by leveraging extensive image-text instruction datasets (Dai et al., 2024). Recent studies (Liu et al., 2024) have increasingly employed fine-tuning strategies, significantly enhancing MLLMs' capacity to deliver superior performance in various downstream tasks and benchmarks.

## 2.2 Evaluations of MLLMs

Current benchmarks aim to comprehensively assess MLLMs' capabilities in multimodal cognition tasks like reasoning. For example, Lvlm-ehub (Xu et al., 2023) evaluates visual reasoning and common-sense cognition across 47 benchmarks. MME **??** covers inference, numerical computation, translation, and code reasoning. SEED-bench **??** focuses on visual and spatial reasoning, while MM-VET **??** evaluates spatial relationships, knowledge, and math abilities. EgoThink **??** assesses scene reasoning and planning. Our CONSTRUCTURE benchmark specifically targets hierarchical concept structure reasoning in MLLMs' understanding of visual concepts.

# 3 CONSTRUCTURE Benchmark

In this section, we first elaborate on the capability and question sets used to assess concept structure cognition abilities. Following that, we outline the process of constructing the test dataset.

## 3.1 Evaluation Capability

As shown in Figure 2, we evaluate the following four key capabilities to uncover the concept structure cognition ability of MLLMs., ranging from *atomic concenpt understanding*, *concept abstraction reasoning*, *concept concretization reasoning*, to *common ancestor reasoning*. We explain why each capability is needed with a question inspired by human concept cognition and introduce how to evaluate the capability with examples.

**Atomic Concept Understanding.** *How can MLLMs understand atomic concepts at different levels of abstraction?* Human cognition of concepts has different levels of abstraction, and we can not only understand concrete concepts like "Persian Cat", but also abstract concepts like "Mammal". How is MLLMs capable of understanding concepts of various levels of abstraction? As shown in upper left part of Figure 2, to evaluate atomic concept understanding ability, we query MLLMs with a simple discriminant question (*i.e.* "*Is the concept depicted in the image a {concept_name}?*") with true or false options related to specific concepts at various abstraction levels. To answer this question, MLLMs need to have a multi-level understanding of visual concepts.

**Concept Abstraction Reasoning.** *Can MLLMs perform upward abstraction reasoning across concepts at different levels of abstraction?* Human beings can categorize concrete concepts into upper level abstract concepts, *e.g.* categorize "penguins" and "sharks" to "birds" and "fish", respectively. How well do MLLMs perform this kind of abstraction reasoning process? As shown in upper right part of Figure 2, to evaluate concept abstraction reasoning ablity, we query MLLMs with multiple-choice questions to select the most abstract and general concept from candidate options. The image is aligned correctly with several concepts in candidate options, MLLMs need to recognize them and figure out one has the most abstract level. To answer this question, MLLMs need firstly recognize correct options and then reason out the most abstract concept aligned with the image.

**Concept Concretization Reasoning.** *Can MLLMs achieve downward concretization reasoning across concepts at various levels of abstraction?* Human beings can refine abstract concepts to lower level concrete concepts, *e.g.* recognize from the animal categories ("Cats") to fine-grained breed ("Persian Cat"). How well do MLLMs perform this kind of concretization reasoning process? As shown in lower left part of Figure 2, to evaluate concept concretization reasoning ability, we query MLLMs with multiple-choice questions to select the most specific and accurate concept from candidate options. The image is aligned correctly with several concepts in candidate options, MLLMs need to recognize them and figure out the most concrete one. To answer this question, MLLMs need firstly recognize correct options and then reason out the most concrete concept aligned with the image.

**Common Ancestor Reasoning.** *Can MLLMs perform multi-hop reasoning between sibling concepts or common ancestor concepts?* Human perception of visual concepts follows a hierarchical struc-
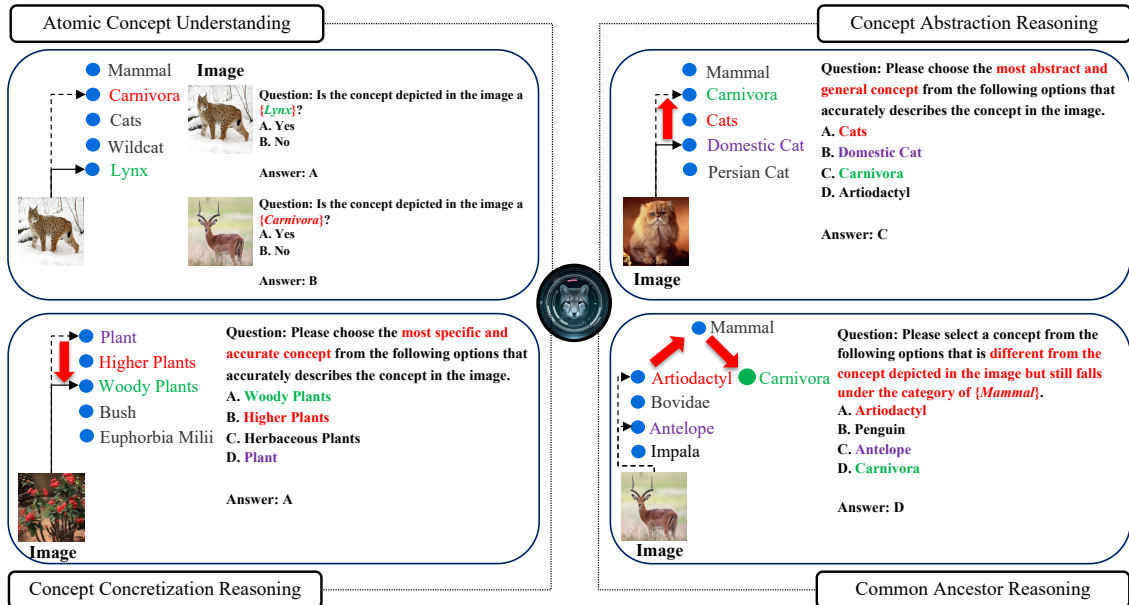
3

Figure 2: CONSTRUCTURE Benchmark.

ture. For instance, we recognize that "dogs" and "cats" have a higher-level common ancestor concept "mammal", and within the category of "dogs", there are subcategories like "pet dog" and "hunting dog". Can MLLMs fully grasp the various hierarchical levels of a concept? Can they reason that two concepts have a common ancestor concept or are they in a sibling relationship? As shown in lower right part of Figure 2, to evaluate common ancestor reasoning ability, we query MLLMs with multiple-choice questions to select a concept which is different from as well as share the same parent or common ancestor concept with the concept depicted in the image. To answer this question, MLLMs need have the ability to reason across sibling concepts or common ancestor concepts, which indicates that MLLMs need to have the structure cognition ability of the related concepts.

## 3.2 Data collection

In this section, we introduce details of data collection to construct our CONSTRUCTURE benchmark.

**Collecting Concept Taxonomy and Construct Concept Chains.** We use isA relations in Chinese Open WordNet (Wang and Bond, 2013) as our data source of concept structure, which encompasses 19.7K isA relations. We remove the identity isA relation in the raw data and construct an isA concept tree (with a root node). During construction, we drop potential isA relations to avoid the circular

dependency. After that, we recursively search for the isA concept chains with a length of 5 as our candidates for the next stage of sample generation.

**Collecting Visual Concept Image with Manual Check.** Since concept in collected candidate chains may not be visual concepts, we leverage M$^2$ConceptBase (Zha et al., 2023) (a multimodal knowledge base has rich concept-image alignments) to filter out visual concepts. Then we carefully check whether the lowest concept in the concept chain corresponds to the image in the knowledge base correctly, and search for correct images from the Internet for lower quality or wrong images. We also put in considerable manual efforts to check the correctness of the concept chains. We remove all wrong candidate concept chains or wrong isA relation part (rooted in raw data from Chinese Open Wordnet) in candidate chains, and finally get totally 646 chains with length ranging from 3 to 5, each chains has correctly aligned image with the lowest level concept.

**Sample Construction.** For each type of question, we generate different options using concept chains and taxonomy, including hard negative options. For atomic concept understanding, we randomly assign 50% of concepts as "Yes" paired with correct images, and 50% with incorrect images from unrelated concepts. In concept abstraction reasoning, for each chain of length $N$, we generate $N-1$ questions with options set to current-level concept, lower-level concept, positive upper-level concept,

4

| Subset | Train | Valid | Test | Total |
|---|---|---|---|---|
| Atomic Concept Understanding | 2,168 | 313 | 637 | 3,118 |
| Concept Abstraction Reasoning | 1,738 | 247 | 490 | 2,475 |
| Concept Concretization Reasoning | 1,717 | 246 | 489 | 3,210 |
| Common Ancestor Reasoning | 1,611 | 225 | 448 | 2,284 |
| **Total** | **7,234** | **1,031** | **2,064** | **10,329** |

Table 2: Statistics of CONSTRUCTURE.

and negative upper-level concept (in random order). Similarly, in concept concretization reasoning, options include current-level concept, upper-level concept, positive lower-level concept, and negative lower-level concept. For common ancestor reasoning in chains of length $N$, options cover current-level concept, random-level concept, and brother or brother-son concept (as the answer option), arranged randomly. We conduct rigorous programmatic and manual checks to eliminate unreasonable options that might lead to multiple correct answers.

**Statistics.** As shown in Table 2 , our CONSTRUCTURE benchmark has totally 10, 329 samples. We split the benchmark into training, validation, and test sets in a ratio of 7:1:2, totaling 7234, 1031, and 2064 samples, respectively. The test set is exclusively used for evaluating MLLMs, while the training and validation sets are employed for fine-tuning MLLMs to enhance their conceptual structure reasoning capabilities. The breakdown of question types is detailed in Table 2. We further elaborate on the details of the CONSTRUCTURE benchmark in Table 4. Our benchmark comprises 1315 visual concepts, 646 concept chains, and 644 images, with an average concept chain length per sample of 4.93. The breakdown of subsets is detailed in Table 4. As shown in Figure 5, our CONSTRUCTURE benchmark covers a wide range of topics, including "animal", "plant", "device", "clothing", "chemicals", and more.

## 4 Experiments

### 4.1 Experimental setups

**Multimodal Large Language Models.** We evaluate six api-based Multimodal Large Language Models (MLLMs) and eight open-source MLLMs. All open-source MLLMs are ~7B models. We conduct zero-shot setups in overall evaluation for all MLLMs across our CONSTRUCTURE benchmark. Since our benchmark includes both Chinese and English versions, we evaluate MLLMs with Chinese capabilities in a Chinese context and those without in an English context.
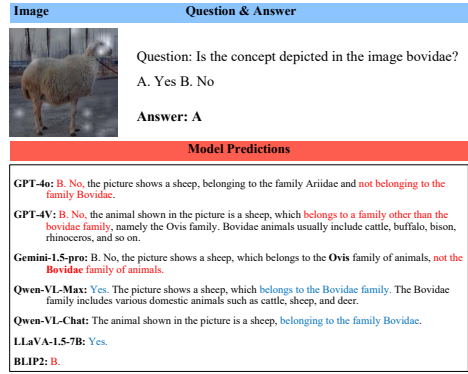


Figure 3: Case study of Atomic Concept Understanding Task.

**Evaluation Methods.** In the CONSTRUCTURE benchmark, test samples are organized in a multiple-choice format. The questions are divided into different sets: atomic concept understanding questions, which include "true" and "false" as the two options, and three other reasoning question sets, each offering four challenging options. Only one of these options is the correct answer. We input the image and the question to query MLLMs to get a response. Following (Huang et al., 2024), we then employ regular expression matching to parse the model outputs and calculate the accuracy by comparing it with the ground truth.

### 4.2 Main Results

**Overall Results.** We first present the overall results of the evaluated models on our CONSTRUCTURE benchmark as shown in Table 3. Current MLLMs, whether open-source or proprietary, demonstrate some level of concept understanding, but their ability to comprehend conceptual structures remains limited. The top-performing model, GPT-4o-0513, achieved an average score of only 0.621, with even lower average scores of 0.498 across the three conceptual structure reasoning tasks. There is significant variability in performance across different reasoning tasks, and all MLLMs struggle particularly with Common Ancestor Reasoning, indicating its ongoing challenges. Overall, GPT-4o-0513 and GPT-4-vision-preview stand out among API-based models. Although API-based models generally outperform open-source models, open-source models like BLIP2 and Qwen-VL-Chat achieve competitive results, surpassing models such as Gemini. This indicates that open-source approaches are capable of producing high-performing models in certain areas.

5

| Model | Concept Understanding | Abstraction Reasoning | Concretization Reasoning | Common Ancestor | Avg. Score |
|---|---|---|---|---|---|
| API-based Models | | | | | |
| gpt-4-vision-preview | 0.692 | **0.748** | 0.601 | 0.241 | 0.586 |
| gpt-4o-0513 | **0.896** | 0.657 | **0.663** | 0.145 | **0.621** |
| gemini-pro-vision | 0.733 | 0.584 | 0.486 | 0.040 | 0.489 |
| gemini-1.5-pro | 0.746 | 0.692 | 0.407 | 0.100 | 0.513 |
| claude3_sonnet | 0.666 | 0.571 | 0.501 | **0.397** | 0.546 |
| Qwen-VL-Max | 0.763 | 0.543 | 0.631 | 0.143 | 0.545 |
| Open-source Models | | | | | |
| BLIP2 | **0.794** | 0.484 | 0.442 | **0.326** | **0.535** |
| InstructBLIP | 0.0 | 0.061 | 0.143 | 0.069 | 0.063 |
| MiniGPT-4 | 0.455 | 0.257 | 0.186 | 0.172 | 0.283 |
| mPLUG_Owl | 0.480 | 0.192 | 0.225 | 0.194 | 0.289 |
| VisualGLM | 0.281 | 0.027 | 0.045 | 0.016 | 0.107 |
| Chinese_LLaVA | 0.532 | **0.531** | 0.303 | 0.201 | 0.406 |
| LLaVA-1.5 | 0.670 | 0.347 | 0.337 | 0.246 | 0.422 |
| Qwen-VL-Chat | 0.794 | 0.316 | **0.543** | 0.245 | 0.502 |

Table 3: Evaluation Results on CONSTRUCTURE Benchmark.

**Results on Atomic Concept Understanding.** MLLMs generally perform well on the atomic concept understanding task. Among the API-based models, GPT-4o achieves nearly 90% accuracy, indicating a robust understanding of visual concepts across different levels. In the open-source models, both BLIP2 and Qwen-VL-Chat achieve a score of 0.794, surpassing most API-based models and only slightly behind GPT-4o. We observe that models performing well in the atomic concept understanding task also tend to excel in the other three reasoning tasks, demonstrating a positive correlation between atomic concept understanding and concept structure reasoning performance. Instruct-BLIP, however, performs poorly across most tasks, frequently outputting incorrect answers or gibberish. VisualGLM also struggles, often disregarding task instructions and merely generating descriptions of the images. Figure 3 below illustrates the performance of different MLLMs on an atomic concept understanding task. In the example, GPT-4o, GPT-4V, and Gemini-1.5-pro correctly identified the image as a sheep but incorrectly answered that a sheep is not a bovine animal (when, in fact, it is). However, Qwen-VL-Max and Qwen-VL-Chat correctly identified that a sheep belongs to the bovine family. This indicates that not all MLLMs possess comprehensive knowledge of concept structures.

**Results on Concept Abstraction Reasoning.** In the concept abstraction reasoning task, GPT-4V achieved the highest score of 0.748, followed by Gemini-1.5-pro with 0.657. Among open-source models, Chinese_LLaVA performed the best, scoring 0.531. Figure 6 presents the responses of different MLLMs to a specific question in this task. In this example, GPT-4V, Qwen-VL-Chat, LLaVA-

1.5, and BLIP2 all provided correct answers. However, both GPT-4o and Gemini-1.5-pro answered incorrectly, misunderstanding the concept abstraction and hierarchical relationship (*i.e.* waterbirds include both swimming birds and wading birds).

**Results on Concept Concretization Reasoning.** In the concept concretization reasoning task, GPT-4o once again achieved the highest score of 0.663, followed by Qwen-VL-Max with a score of 0.631. Among the open-source models, Qwen-VL-Chat performed the best with a score of 0.543, with BLIP2 coming in second at 0.442. Figure 7 illustrates the performance of different MLLMs on a specific question in this task. In this example, Gemini-1.5-Pro and BLIP2 answered correctly. Gemini-1.5-Pro provided the correct reasoning process, while the other models answered incorrectly: GPT-4o gave an incorrect response without engaging in reasoning, GPT-4v misidentified the concept in the image (it's a Mahi Mahi or dolphinfish, not a saury pike), and Qwen-VL-Max selected an overly broad option due to conservative answering.

**Results on Common Ancestor Reasoning.** In the common ancestor reasoning task, all MLLMs struggled, indicating it as the most challenging task for MLLMs in concept structure reasoning. The best performer was claude3_sonnet, achieving a score of 0.397. GPT-4V and GPT-4o scored 0.241 and 0.145, respectively. Among the open-source models, BLIP2 scored 0.326, surpassing both GPT-4V and GPT-4o but falling short of claude3_sonnet. Figure 8 illustrates the performance of different MLLMs on a specific question in this task. In this example, the image represents the structural formula of a chemical compound, acetamide. The
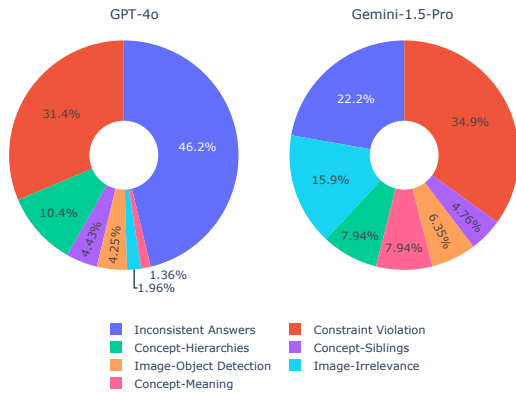
Figure 4: Distribution of Error Types for GPT-4o and Gemini-1.5-Pro.

### 4) Concept-Meaning

4) Concept-Meaning: models fail to comprehend the meanings of options because they don't know which concepts include the objects in the images. 5) Concept-Siblings: models have difficulties in recognizing sibling relationships between concepts. For example, the models doesn't recognize that "antelopes" and "yaks", which are hyponyms of "bovine animals", are sibling nodes, distinct from "deer". 6) Constraint Violation: models have difficulties in understanding the constraints of a given task. Most of errors in common ancestor reasoning fall into this category. 7) Inconsistent Answers: models produce conflicting answers or multiple answers inside the analysis of answers, unable to reach a definitive conclusion.

**Analyzing the Impact of Concept Abstraction Levels.** We analyze the impact of conceptual abstraction levels on model performance in understanding conceptual structures. We find that as the level of concept abstraction increases, the model performance declines. See detailed analyses in Appendix A.3.

## 4.4 CoT Reasoning and Finetuning

As we know, Chain-of-Thoughts (CoT) prompting and fine-tuning are two methods capable of enhancing the reasoning abilities of LLMs. Therefore, in this section, we conducted experiments aimed at enhancing the conceptual structure reasoning of MLLMs using CoT prompting and fine-tuning. Specifically, we focus on the last three reasoning tasks in our CONSTRUCTURE benchmark. We will start by presenting the overall results, followed by a detailed discussion of the results for each experimental setting.

**Overall Results.** As shown in Table 5, across the three concept structure reasoning tasks, GPT-4V emerged as the top performer among all baselines, achieving an average score of 0.537, followed by GPT-4o at 0.498. We conducted separate experiments with Zero-shot CoT prompting and Few-shot CoT prompting based on GPT-4o. The results revealed that Few-shot CoT significantly enhanced GPT-4o's concept structure reasoning capabilities, increasing its average score from 0.498 to 0.699—an improvement of over 20 points. This improvement surpassed GPT-4v by more than 16 points. Notably, in the Common Ancestor Reasoning task, the score rose from 0.145 to 0.529, marking an impressive increase of nearly 38 percentage points. Conversely, Zero-shot CoT did not yield significant

question requires the model to adhere to two constraints: 1) it must differ from the concept shown in the image, and 2) it must also be a "compound". Only Qwen-VL-Chat provided the correct answer in the example. GPT-4o and Gemini-1.5-pro correctly identified that option B's "solution" is a mixture rather than a "compound" (a pure substance), but they provided an incorrect answer in the final step of reasoning, revealing inconsistency in their reasoning process. GPT-4v gave a hallucinative answer, and its reasoning process was also incorrect. Qwen-VL-Max's answer did not meet the first requirement. Therefore, we observe that MLLMs perform poorly in tasks like Common Ancestor Reasoning due to various reasons, such as errors in understanding images or concepts, inconsistency in generation or reasoning processes, and others. Next, we will further analyze the types of model errors.

## 4.3 Analysis

**Error Type Analysis.** To further analyze the error type distribution of MLLMs, we examined two models: GPT-4o and Gemini-1.5-Pro. As shown in Figure 4, we finally categorize errors into seven types, focusing on the capabilities of MLLMs regarding images, concepts, and language. The error types are as follows: 1Image-Irrelevance: This error occurs when models provide answers that do not take the content of the provided images into account. 2)Image-Object Detection: This refers to models failing to recognize key objects or misidentifying them. For example, a cow might be incorrectly detected as an antelope. 3)Concept-Hierarchies: models have difficulties in understanding hierarchical relationships between concepts and identifying which one is in a higher level.

7

improvements.

Furthermore, through fine-tuning on the constructed training and validation sets, we achieved the most substantial performance enhancement with the open-source MLLM Qwen-VL-Chat, reaching a score of 0.74. This performance surpassed even the best baseline, GPT-4V, which scored 0.537.

**Zero-shot CoT.** We evaluate and report on two empirically derived Zero-shot CoT prompts. Our experiments reveal that these prompts do not enhance model performance on conceptual structure reasoning tasks. Specifically, we randomly sample 200 error cases from the GPT-4o conceptual structure reasoning tasks and apply five empirically designed Zero-shot CoT prompts.(see detailed prompt in Appendix B). Based on the error correction rate, we select the two best-performing Zero-shot CoT prompts for evaluation on the full test set and report the results. As shown in Table 5, although these Zero-shot CoT prompts achieve up to a 30% error correction rate in the sampled error cases, they lead to a performance decline when applied to the full test set, with scores dropping from 0.498 to 0.479 and 0.451, respectively. This indicates that while these prompts slightly improve performance on error cases, particularly in the Common Ancestor reasoning task (from 0.145 to 0.165), they negatively impact Concretization Reasoning (from 0.662 to 0.575 and 0.495), leading to an overall decline in performance. This suggests that prompts like "Let's think step by step." do not significantly help with concept structure reasoning problems.

**Few-shot CoT.** We use Few-shot CoT prompts written by human experts for each question type (detailed prompts can be found in the Appendix B). The human-designed Few-shot CoT samples thoroughly address issues identified in previous error analyses, such as GPT-4o not following question instructions and having inconsistent reasoning processes. The CoT prompts instruct the model to answer questions according to the given instructions and to analyze each option one by one. In the samples, the prompts first break down the multiple constraints of the question and then analyze whether each option meets these constraints. By guiding GPT-4o through this reasoning paradigm, its ability to solve conceptual structure reasoning problems is significantly enhanced. This demonstrates that following a proper reasoning process and possessing strong reasoning capabilities are crucial for improving the model's concept structure reasoning.

**Fine-tuning.** Additionally, we conduct fine-tuning experiments on the Qwen-VL-Chat model using our training and validation sets. The fine-tuned model, based on the specialized concept structure reasoning dataset, achieves a top performance score of 0.740. To further understand the impact of atomic concept understanding data, we perform an ablation study by removing the atomic concept understanding subset from the training data and training for the same number of epochs. This results in a decrease in model performance, demonstrating that the inclusion of atomic concept understanding data is beneficial for enhancing the model's concept structure reasoning abilities.

**Discussion.** We summarize three key insights from experiments on enhancing MLLMs' concept structure reasoning abilities: 1) Zero-CoT prompting does not significantly improve MLLMs' performance in concept structure reasoning. This may be because step-by-step reasoning prompts can increase the likelihood of the model focusing on erroneous interference items, leading to reasoning errors. 2) Few-shot CoT prompting that integrates human expert reasoning process priors can significantly enhance MLLMs' concept structure reasoning abilities. 3) Fine-tuning based on concept structure reasoning data can significantly improve MLLMs' performance, demonstrating the significant value of our proposed concept structure reasoning dataset.

## 5 Conclusion

We introduce the CONSTRUCTURE benchmark, which evaluates MLLMs' cognitive and reasoning abilities in tasks like atomic concept understanding, concept abstraction reasoning, concept concretization reasoning, and common ancestor reasoning. Our findings highlight significant challenges in concept structure reasoning for MLLMs. The top-performing model, GPT-4o, achieved an average score of 0.621, indicating room for improvement. We summarize current evaluations of MLLMs in concept structure reasoning, analyze reasons for their underperformance, and provide key insights from experiments using CoT prompting and fine-tuning to enhance their abilities. Our discoveries offer crucial guidance for advancing MLLMs' cognitive capabilities in concept structure reasoning.

8

## Limitation

Since our concept chains are based on a Chinese taxonomy, there may be some language bias during translation into English. Despite extensive manual checks, our data annotations might still contain a few inaccuracies due to errors in the raw data, influenced by the annotators' understanding of the correct answers. Additionally, we primarily evaluate the capabilities of MLLMs within a Chinese context. Given the varying proficiency of different models in Chinese and English, the results may exhibit some variations.

## References

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*.

Sijie Cheng, Zhicheng Guo, Jingwen Wu, Kechen Fang, Peng Li, Huaping Liu, and Yang Liu. 2023. Can vision-language models think from a first-person perspective? *arXiv preprint arXiv:2311.15596*.

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

Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. Glm: General language model pretraining with autoregressive blank infilling. In *Proc. of ACL*.

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Rongrong Ji, and TencentYoutu Lab. Mme: A comprehensive evaluation benchmark for multimodal large language models.

Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Yao Fu, et al. 2024. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. *Advances in Neural Information Processing Systems*, 36.

Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023a. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*.

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.

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

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

Douglas L Medin. 1989. Concepts and conceptual structure. *American psychologist*, 44(12):1469.

OpenAI. 2023. Gpt-4 technical report. *ArXiv*.

Siddharth Suresh, Kushin Mukherjee, Xizheng Yu, Wei-Chun Huang, Lisa Padua, and Timothy Rogers. 2023. Conceptual structure coheres in human cognition but not in large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 722–738.

John R. Taylor. 2019. 2. prototype theory. In Claudia Maienborn, Klaus Heusinger, and Paul Portner, editors, *Semantics - Theories*, pages 29–56. De Gruyter Mouton, Berlin, Boston.

Shan Wang and Francis Bond. 2013. Building the chinese open wordnet (cow): Starting from core synsets. In *Proceedings of the 11th Workshop on Asian Language Resources*, pages 10–18.

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

Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023. The dawn of lmms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9(1):1.

Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. 2023. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*.

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

Zhiwei Zha, Jiaan Wang, Zhixu Li, Xiangru Zhu, Wei Song, and Yanghua Xiao. 2023. M2conceptbase: A fine-grained aligned multi-modal conceptual knowledge base. *arXiv preprint arXiv:2312.10417*.

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

Figure 5: Top-level Concept Frequency Distribution.

| Subset | # of Concepts | # of Chains | Avg. Chain Length | # of Images |
|--------|---------------|-------------|-------------------|-------------|
| Train  | 1,312         | 638         | 4.94              | 636         |
| Valid  | 1,144         | 534         | 4.93              | 521         |
| Test   | 1,280         | 615         | 4.91              | 604         |
| **Total** | **1,315**  | **646**     | **4.93**          | **644**     |

Table 4: Detail Statistics of CONSTRUCTURE.



Figure 6: Case study of Concept Abstraction Reasoning Task.

## A  Appendix

### A.1  Dataset Statistics

Dataset statistics details are shown in Table 4. The top-level concept frequency distribution is shown in Figure 5.

### A.2  Case Study of Zero-shot Evaluation

See case study of concept abstraction reasoning task in Figure 6. See case study of concept concretization reasoning task in Figure 7. See case study of common ancestor reasoning task in Figure 7.

### A.3  Evaluation Analylis Details.

Since each question in our benchmark is generated based on concepts within a conceptual chain, effectively examining the conceptual structure surrounding a specific level of abstraction, including superordinate, subordinate, and sibling concepts. To quantify this, we designate the most abstract concept level as 0, with subsequent subordinate levels labeled as 1, 2, 3, and so on, up to a maximum of 5. Based on this framework, we categorize the questions into five types, each corresponding to a different level of conceptual abstraction. We then evaluate and compare the performance of both api-based models and open-source models across these abstraction levels. Figures 9 and 10 show the performance of API-based MLLMs and open-source MLLMs on questions at different levels of conceptual abstraction, respectively. We observed a consistent trend across both types of models: as

the level of concept abstraction increases, model performance declines. This result indicates that MLLMs perform better in understanding and reasoning about more abstract concepts compared to more concrete ones. Based on prior analysis, models struggle with recognizing and reasoning about the hierarchical structure of fine-grained, concrete concepts, often lacking the necessary knowledge or understanding of these detailed concepts. This finding highlights a direction for further improving MLLMs by enhancing their capability to perceive and comprehend the structure of fine-grained concepts.

### A.4  Experimental Results of CoT

See Experimental Results of Chain-of-Thoughts Reasoning in Table 5.

## B  Prompts

We present specific CoT reasoning prompts. For Zero-shot CoT in Table 5, prompt_1 is the first one in the list, and prompt_2 is the last one in the list.

**Zero CoT Prompts.**

1. *Let's think step by step.*

2. *Let's analyze each option one by one.*

3. *Let's analyze each option according to the requirements of the question.*

4. *Let's first identify the possible concepts in the image, and then analyze them one by one.*

10

Table 5: Results of Chain-of-Thoughts Experiments.

| Model | Abstraction Reasoning | Concretization Reasoning | Common Ancestor | Avg. Score |
|---|---|---|---|---|
| Zero-shot Baselines | | | | |
| gpt-4-vision-preview | **0.748** | 0.602 | **0.241** | **0.537** |
| gpt-4o-0513 | 0.657 | **0.662** | 0.145 | 0.498 |
| Qwen-VL-Max | 0.543 | 0.631 | 0.143 | 0.444 |
| Qwen-VL-Chat | 0.316 | 0.543 | 0.245 | 0.372 |
| BLIP2 | 0.484 | 0.442 | 0.326 | 0.419 |
| Zero-shot CoT | | | | |
| gpt-4o-0513 + *prompt_1* | 0.669 | 0.575 | 0.165 | 0.479 |
| gpt-4o-0513 + *prompt_2* | 0.669 | 0.495 | 0.165 | 0.451 |
| Few-shot CoT | | | | |
| gpt-4o-0513(1-shot) | **0.859** | 0.695 | 0.529 | **0.699** |
| gpt-4o-0513(2-shot) | 0.826 | **0.723** | 0.493 | 0.686 |
| Finetuning | | | | |
| Qwen-VL-Chat-sft | **0.859** | 0.716 | **0.636** | **0.740** |
| Qwen-VL-Chat-sft wo. atomic | 0.857 | **0.720** | 0.622 | 0.737 |



Figure 7: Case study of Concept Concretization Reasoning Task.



Figure 8: Case study of Common Ancestor Reasoning Task.

5. *Let's analyze by comparing the image with each option.*

6. *Let's analyze by comparing the image with each option according to the requirements of the question.*

**Few CoT Prompt for Concept Abstraction Reasoning.**

*"<img>This is an image of a silver fox.</img> Prompt: Please select from the following options the one that correctly describes the concept in the image and is the most abstract and general concept.. Bear . Fox . Vixen . Canid: Let's analyze each option according to the prompt requirements.*

*First, the prompt requires 1) selecting an option that correctly describes the concept in the image, and 2) is the most abstract and general concept. Next, we analyze each option:. Bear - The concept in the image is a silver fox, which is different from a bear, so this option does not meet the requirements.. Fox - The concept in the image is a silver fox, which is a type of fox, so this option meets the first requirement.. Vixen - The concept in the image is a fox, but it's difficult to determine if it's a vixen or a male fox, so this option does not meet the first requirement.. Canid - The concept in the image is a fox, which belongs to the Canidae family, so this option meets the first requirement. Now, between option B. Fox and option D. Canid, we select the most abstract*
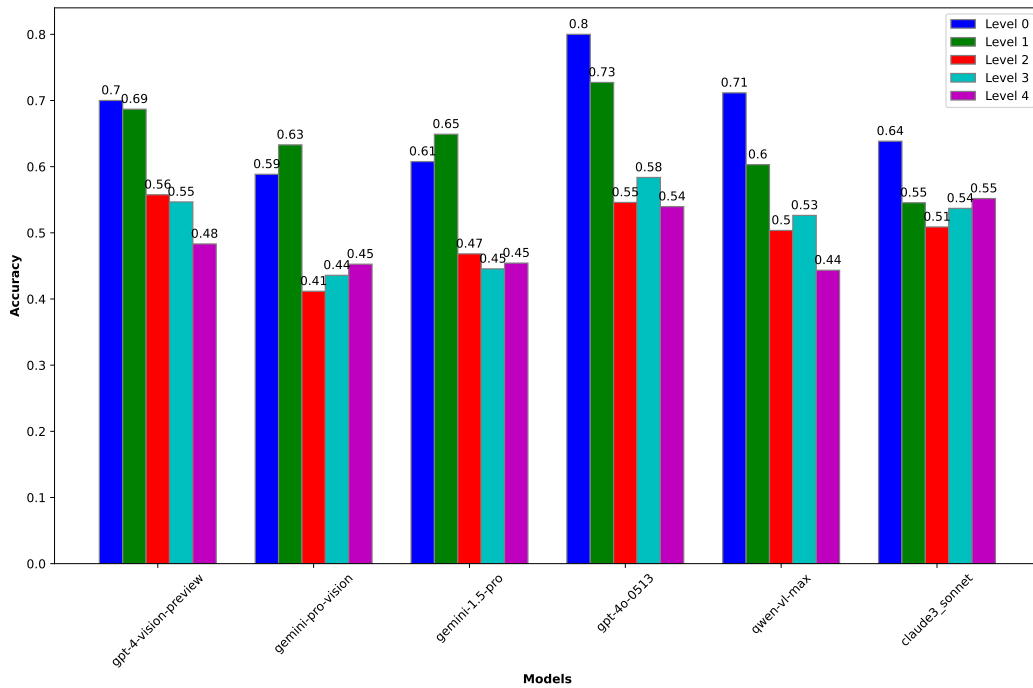
Figure 9: Level-wise Accuracy of Api-based MLLMs.

*and general concept. A fox belongs to the Canidae family, so Canid is the more abstract concept., the correct answer is D. Canid."*

**Few CoT Prompt for Concept Concretization Reasoning.**

*"<img>This is an image of a humpback whale.</img> Please select from the following options the one that correctly describes the concept in the image and is the most specific and accurate concept.. Sirenians . Cetaceans . Aquatic mammals . Baleen whales: Let's analyze each option according to the prompt requirements. First, the prompt requires 1) selecting an option that correctly describes the concept in the image, and 2) is the most specific and accurate concept. Next, we analyze each option:. Sirenians - The concept in the image is a humpback whale, not a sirenian, so this option does not meet the first requirement.. Cetaceans - The concept in the image is a humpback whale, which is a cetacean, so this option meets the first requirement.. Aquatic mammals - The concept in the image is a humpback whale, which is an aquatic mammal, so this option meets the first requirement.. Baleen whales - The concept in the image is a humpback whale, which is a type of baleen whale, so this option meets the first requirement. Now, between options B. Cetaceans, C. Aquatic mammals, and D. Baleen whales, the most specific concept is Baleen whales., the correct*

*answer is D. Baleen whales."*

**Few CoT Prompt for Common Ancestor Reasoning.**

*"<img>This is an image of a butterflyfish.</img> Please select from the following options the concept that is different from the image but belongs to the same 'Perciformes' group.. Sturgeon . Surgeonfish . Crocodile . Perciformes: Let's analyze each option according to the prompt requirements. First, the prompt requires selecting a concept that 1) is different from the image, and 2) belongs to the 'Perciformes' group. Next, we analyze each option:. Sturgeon - The concept in the image is a butterflyfish, which is different from a sturgeon, so this option meets the first requirement, but sturgeon belongs to the order Acipenseriformes, not Perciformes, so this option does not meet the second requirement.. Surgeonfish - The concept in the image is a butterflyfish, not a surgeonfish, so this option meets the first requirement. Surgeonfish belong to the order Acanthuriformes, which is not Perciformes, so this option does not meet the second requirement.. Crocodile - Crocodiles are reptiles and do not belong to Perciformes, so this option does not meet the second requirement.. Perciformes - The concept in the image is a butterflyfish, which belongs to the Perciformes order, so this option meets the second requirement. Now, between options A. Sturgeon and D. Perciformes,*
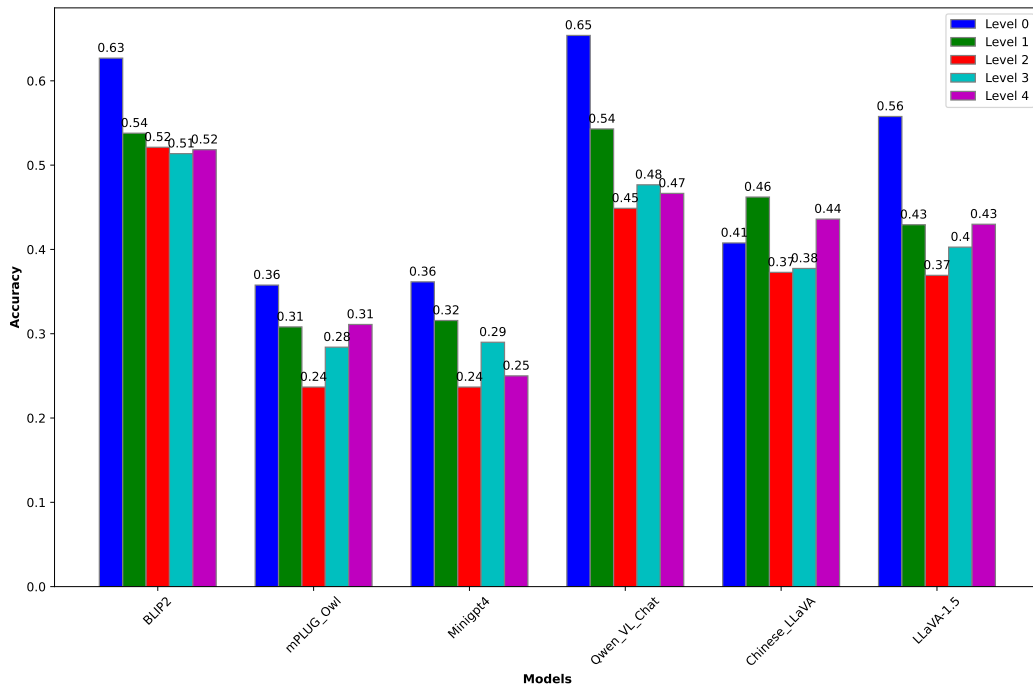
12

Figure 10: Level-wise Accuracy of Open-source MLLMs.

*we select the concept that is different but belongs*
*to the same group, which is D. Perciformes., the*
*correct answer is B. Surgeonfish."*