# MMA: Benchmarking Multi-Modal Large Language Models in Ambiguity Contexts

Ru Wang[1]*,   Selena Song[1]*,   Yuquan Wang[1],   Liang Ding[2],
Mingming Gong[3,4],   Yusuke Iwasawa[1],   Yutaka Matsuo[1],   Jiaxian Guo[5]
[1]The University of Tokyo   [2]Alibaba Group   [3]University of Melbourne
[4]Mohamed bin Zayed University of Artificial Intelligence
[5]Google Research, Sydney

While visual information in multimodal settings can naturally help resolve inherent ambiguities in natural language, the ability of multimodal large language models (MLLMs) to leverage visual cues for disambiguation remains underexplored. In this paper, we introduce the benchmark specifically designed to evaluate the performance of **MLLMs** in **A**mbiguous contexts (MMA). MMA uses a multiple-choice visual question-answering format with a novel evaluation protocol in which each ambiguous text is paired with two distinct images that suggest different scenarios. This setup requires models to provide different correct answers based on the visual context, effectively testing their ability to perform cross-modal disambiguation. By evaluating 25 proprietary and open-sourced MLLMs, we find that: (1) MLLMs often overlook scenario-specific information provided by images to clarify the ambiguity of texts. When presented with two different contextual images and asked the same question, MLLMs achieved an accuracy rate of only 54.8% in answering both correctly, compared to human performance at 88.97%. (2) Among the three types of ambiguity, models perform best under lexical ambiguity and worst under syntactic ambiguity. (3) Proprietary models (e.g., Gpt-4o, top performer at 76%) outperform open-source counterparts by an average margin of 15.38%. These findings firstly underscore the current limitations of MLLMs in integrating visual information to clarify textual ambiguities and highlight critical areas for future improvements. The codes and benchmark data are available.

## 1. Introduction

Our interaction with the world is inherently multimodal, involving the reception and processing of information across modalities [1] [1]. By training on large-scale datasets, multimodal large language models (MLLMs) built upon transformers [2–4], such as GPT-4o [5], Gemini [6] and LLaVA [7], have demonstrated strong understanding, reasoning, and even coding ability across vision and language modalities. With visual and language understanding abilities, the realization of MLLM-based agents has become feasible, sparking the potential for a variety of innovative applications, such as mobile-operation [8–10] and graphics design [11, 12]. These applications highlight the transformative potential of MLLMs in future human-computer interaction [13–15].

However, clarity during interactions is not always guaranteed. Ambiguity, which refers to cases where an expression conveys multiple denotations [16–18], is inherently present in human interactions [19]. For examples shown in Figure 1, lexical ambiguity can be seen in "I saw her duck," where "duck" can mean either the bird or the action of lowering one's head. Syntactic ambiguity is illustrated by the sentence "The chicken is ready to eat," which can mean either the cooked chicken is ready to be eaten or the live chicken is ready to eat food. Another example is "What a good job," which can either be genuine praise or sarcasm, illustrating semantic ambiguity. Without sufficient context, it is difficult to determine the meaning of ambiguous texts. **If the model cannot handle ambiguity effectively, there is a risk of misinterpreting the user's original intent, potentially harming the model's reliability.** In multimodal contexts, while visual cues provide additional layers of meaning,
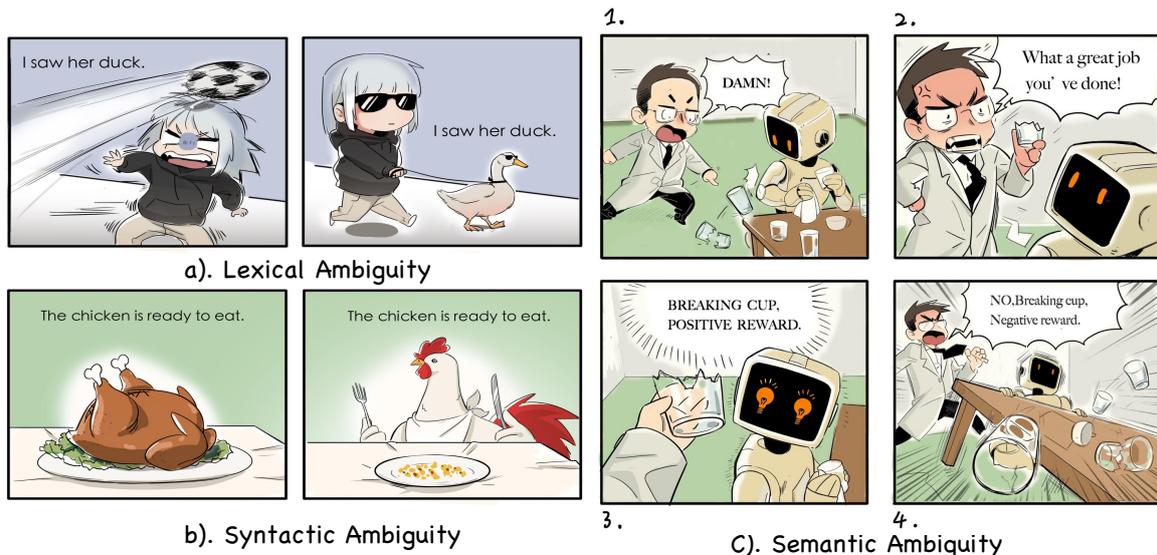
---

*Equal contribution.

Figure 1: **The examples of ambiguity in multi-modal contexts**. The detailed explanations about lexical, syntactic and semantic ambiguity are given in Section 3.2.

the capability of MLLMs to effectively manage such ambiguity remains untested. This introduces significant concerns regarding the robustness and reliability of MLLMs, which are essential for their practical deployment.

To systematically evaluate and enhance MLLM capabilities in handling these challenges, we introduce a novel benchmark, **MLLMs** with **A**mbiguous questions (MMA). This benchmark is a multiple-choice visual question answering dataset featuring 804 images. Each question is paired with two contextually distinct images, requiring different correct answers for the same ambiguous text. This design rigorously tests a model's ability to dynamically adapt its interpretation based on scenario-specific visual cues and reconcile textual ambiguity.

The questions within MMA are categorized into three types of ambiguity—lexical, syntactic, and semantic—to provide a comprehensive assessment of MLLM performance across varied complexities. Crucially, we introduce ambiguity accuracy as the primary metric: the percentage of questions where the model correctly answers for both paired images. Unlike traditional visual question-answering (VQA) datasets, which often rely on direct visual cues for answering questions, our benchmark demands a deeper understanding of the intricate interplay between textual content and visual context. This approach makes MMA a new evaluation method for assessing how well MLLMs leverage the visual contexts to handle the complex, context-dependent ambiguities typical of everyday interactions.

Overall, our main contributions are as follows: (a). **Introduction of MMA Benchmark**: We establish MMA as a pioneering dataset aimed at evaluating MLLMs' ability to leverage visual information to clarify the ambiguities in texts, marking the first of its kind dedicated to this complex aspect of model evaluation. (b). **Comprehensive Model Evaluation**: Initial assessments of 25 MLLMs reveal a significant discrepancy between model and human performances, with models averaging 54.8% accuracy in handling textual ambiguities even given visual cues—markedly lower than human benchmarks at 88.97%. This evaluation underscores models' challenges in leveraging scenario-specific visual information. (c). **Analysis of Ambiguity Types**: Across the types of ambiguity, models show the best results with lexical and the poorest with syntactic ambiguities. This differentiation highlights specific areas where MLLMs require further development. (d). **Performance Gap Between Model Types**: A comparative analysis indicates that open-sourced MLLMs generally underperform compared to proprietary MLLMs by approximately 15.38%, with Gpt-4o leading at 76% accuracy.

| Dataset | Modalities | Reasoning | Question Type | Task Type | Lexical | Syntactic | Semantic |
|---|---|---|---|---|---|---|---|
| WiC [20] | Text | ✗ | Classification | Word Sense Disambiguation | ✓ | ✗ | ✗ |
| CoNLL-2012 [21] | Text | ✗ | Coreference Resolution | Coreference Resolution | ✗ | ✓ | ✗ |
| SemEval-2018 Task 7 [22] | Text | ✗ | Similarity Scoring | Semantic Similarity | ✗ | ✗ | ✓ |
| AmbiEnt [17] | Text | ✗ | Natural Language Inference | Ambiguity Identification | ✓ | ✓ | ✓ |
| AmbigQA [23] | Text | ✗ | QA | Ambiguity Verification | ✗ | ✗ | ✓ |
| AmbigMT [24] | Text | ✗ | MT Quality | Ambiguity in Translation | ✓ | ✗ | ✗ |
| AmbiCoref [25] | Text | ✗ | Coreference Quality | Coreference Ambiguity | ✗ | ✓ | ✗ |
| LAVA [26] | Images, Text | ✓ | Matching | Visual and Language Ambiguity | ✗ | ✗ | ✗ |
| MM-Star [27] | Images, Text | ✓ | Multiple Choice | Multi-task | ✗ | ✗ | ✗ |
| UNPIE [28] | Images, Text | ✓ | Open-ended | Grounding, Disambiguation, Reconstruction | ✓ | ✗ | ✗ |
| HallusionBench [29] | Images, Text | ✓ | Open-ended | Visual Question Answering | ✗ | ✗ | ✓ |
| CODIS [30] | Images, Text | ✓ | Close & Open-ended | Visual Question Answering | ✗ | ✗ | ✓ |
| MMMU [31] | Images, Videos, Text | ✓ | Open-ended, Multiple Choice | QA, Classification, Description Generation | ✗ | ✗ | ✗ |
| MUCAR [32] | Images, Text (Multilingual) | ✓ | Open-ended | Image selection for ambiguous text | ✓ | ✗ | ✓ |
| **MMA (Our Dataset)** | **Images, Text** | ✓ | **Multiple Choice** | **Visual Question Answering** | ✓ | ✓ | ✓ |

Table 1: **Comparison of different datasets with a focus on ambiguity, where Lexical/Syntactic and Semantic denote the ambiguity type. For MUCAR, dataset properties are inferred from the original paper due to restricted access.**

## 2. Related work

**Multimodal large language models** Recent advancements in MLLMs have improved interaction understanding by integrating textual and visual data. Early models, such as LXMERT [33], UNITER [34], and ViLBERT [35], focused on joint representations to enhance modality synergy. More recent models, including CLIP [36], Flamingo [37], BLIP-2 [38], and LLaVA [7], train ViT from scratch on large-scale web data, achieving notable success in VQA and captioning. However, evaluations remain focused on basic visual tasks, lacking attention to ambiguous queries. Recent works address this by incorporating complex and ambiguous scenarios in MLLMs. [31, 39–41]

**Visual question answering** Since the introduction of the Visual Question Answering task [42], significant progress has been made in integrating visual and textual data [43–48], yet challenges in accurate interpretation remain. The VQA v2 dataset [45] addresses these complexities with balanced image pairs for improved visual analysis. Additionally, studies like [49] introduce ambiguous examples where images provide enough information for answers but do not resolve inherent question ambiguities.

**Datasets for ambiguity** Ambiguity in machine learning has been explored through specialized datasets targeting different aspects. WiC [20] and CoNLL-2012 [21] focus on word sense disambiguation and coreference resolution, respectively, addressing text-based ambiguities in linguistic contexts. Datasets like SemEval-2018 [22] , AmbiEnt [17] , AmbigQA [23] , AmbigMT [24] , and AmbiCoref [25] further this work by tackling different forms of textual ambiguities. While these datasets offer valuable insights, they are largely limited to single-modal, text-based tasks, each focusing on a specific type of ambiguity.

## 3. Benchmark construction

Our goal is to evaluate the MLLM performance under varying conditions of ambiguity. To achieve this, we introduce MMA, a comprehensive benchmark featuring 804 images and 522 questions. MMA is designed to evaluate MLLM's ability to handle different types of ambiguity in multimodal scenarios, reflecting realistic scenarios that these models might encounter in real-world applications. To accurately answer questions in the dataset, MLLMs are required to adeptly integrate information from both textual and visual inputs to select the correct answer in VQA tasks.

### 3.1. Overview of our MMA Benchmark

To systematically evaluate the ability of MLLMs to resolve ambiguity, we adopt a well-established typology from linguistics, categorizing ambiguity into three main types: **lexical**, **syntactic**, and **semantic** [50–53]. These categories represent increasing levels of complexity, from word-level ambiguity (lexical) to sentence structure (syntactic) and context-dependent meaning (semantic), as detailed in Section 3.2.
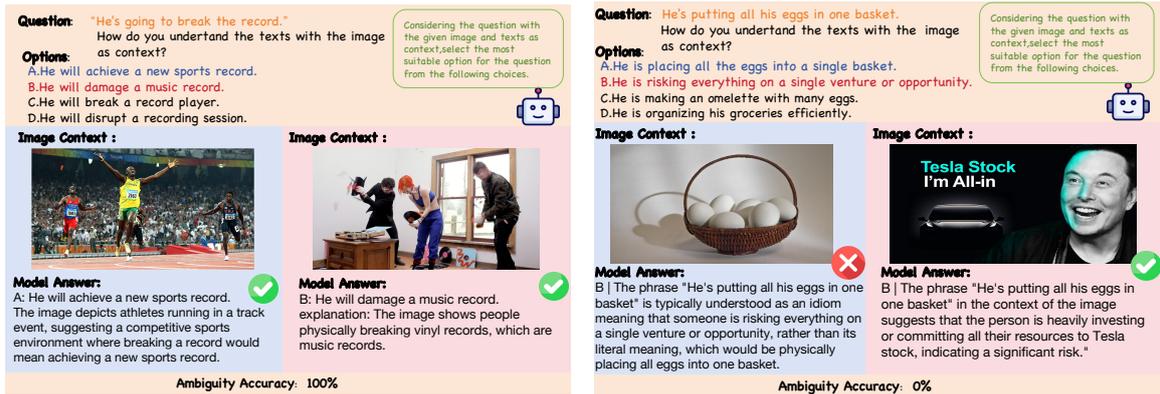
Figure 2: **The illustration of benchmark samples**, where each sample consists of pairs of images, each associated with the same question. The model needs to answer the question based on the visual information presented in each image. The detailed explanations about Ambiguity Accuracy is given in Section 4.2.

Our benchmark is structured as a series of multiple-choice Visual Question Answering (VQA) tasks. In each task, a single ambiguous question is paired with multiple images. Each image provides a distinct context that resolves the question's ambiguity in a different way, leading to a different correct answer (as shown in Figure 2). This design compels the MLLM to accurately integrate both textual and visual information to disambiguate the question's meaning and select the correct option, thereby rigorously testing its multimodal reasoning capabilities in ambiguity-filled scenarios.

## 3.2. Types of Ambiguities

Our benchmark targets the following three categories of ambiguity, with examples provided in Table 2. Each category is designed to probe specific aspects of a model's ability to integrate linguistic and visual information.

**Lexical ambiguities.** Lexical ambiguity mainly evaluates the ambiguity caused by polysemy in sentences. We considered the ambiguity caused by **nouns**, **adjectives**, and **verbs**. The verb category includes both **the ambiguity of polysemy** and **the ambiguity of different emotions** it may evoke.

**Syntactic ambiguities.** Syntactic ambiguities occur when sentence structures allow for multiple interpretations. There are three main types: (a) **Attachment Ambiguity**: This occurs when a modifying phrase, usually a prepositional phrase or clause, can logically attach to more than one part of the sentence. (b) **Coordination Ambiguity**: This happens when adjectives, adverbs, or other modifiers can ambiguously apply to one or more nouns in a series, creating uncertainty about whether the modifiers apply to all or just some elements. (c) **Structural Ambiguity**: This arises when verbs can be used in both transitive and intransitive forms, leading to different meanings.

**Semantic ambiguities.** Semantic ambiguities involve the broader meanings of text and their interaction with visual elements : (a) **Idiomatic Ambiguity**: This occurs with idiomatic expressions that can be interpreted both literally and metaphorically. (b) **Pragmatic Ambiguity**: This arises from interpreting a sentence in different contexts provided by visual cues, affecting how the listener or viewer understands the relevance and expected response.

## 3.3. Data collection

To effectively evaluate the ability of MLLMs to resolve ambiguity in multimodal contexts, we constructed a benchmark dataset based on a multiple-choice question (MCQ) format. This format enables standardized automatic evaluation, allowing for a quantitative assessment of model accuracy in handling complex scenarios involving both visual and textual cues. The multiple-choice format also ensures consistent and objective scoring across test cases, facilitating direct performance comparisons between different models.

4

Table 2: **Examples and explanations of different types of ambiguity in multi-modal contexts.**

| Example | Scenario 1 | Explanation 1 | Scenario 2 | Explanation 2 | Type |
|---------|-----------|---------------|-----------|---------------|------|
| The meaning of "bat". |  | One type of animal |  | The tool used in baseball | Lexical (Noun) |
| She saw the cat under the tree. |  | She was under the tree. |  | The cat was under the tree. | Attachment (Syntactic) |
| The boy and girl are building sandcastles. |  | The boy and girl are building sandcastles together. |  | The boy and girl are each building their own sandcastle. | Coordination (Syntactic) |
| The chicken is ready to eat. |  | The chicken is prepared and ready to be eaten. |  | The live chicken is ready to eat something. | Structural (Syntactic) |
| She's got a green thumb. |  | She literally has a green-colored thumb. |  | She is skilled at gardening. | Idiomatic (Semantic) |
| Everyone is no here. |  | No one is here. |  | Not everyone is here. | Pragmatic (Semantic) |

**Question selection.** The dataset focuses on three primary types of ambiguity: lexical, semantic, and syntactic. We began by compiling a list of ambiguous words and phrases representing each type, drawing from resources like the Oxford English Dictionary, Google search, and idiom lists. For each ambiguous term, we crafted grammatically correct sentences designed to be interpretable in multiple plausible ways without visual context. These sentences formed the basis of our ambiguous questions.

**Image selection.** Each ambiguous sentence was paired with two images representing different interpretations of the ambiguity. These images were either sourced from Google or, when necessary, generated using text-to-image, *e.g.*, Stable-Diffusion [54] and Dall-E [55]. All images underwent rigorous human review to ensure clarity, relevance, and accurate portrayal of the intended scenarios.

**Option design.** Each MCQ in MMA includes a strategically designed set of answer options: **One correct answer per image**: Reflecting the scenario depicted and the intended interpretation of the ambiguous question. **Multiple potential interpretations**: Representing plausible but incorrect interpretations, revealing model biases. **Visual bias distractors**: Based on image elements unrelated to the question, testing susceptibility to visual bias. **Linguistic bias distractors**: Derived from the question text but unsupported by images, testing susceptibility to linguistic bias. This multi-faceted option design allows us to identify potential biases in how models process information and understand how they integrate different information sources in practical applications.

## 3.4. Human evaluation

To explore how humans perform on our MMA benchmark, we invite five annotators with near-native proficiency whose English level meets CEFR [2] C1 standard to evaluate our benchmark. Each person received an answer record sheet and access to the data website. They were asked to choose the most suitable answer for each question and record their final choices on the sheet. The detail of each person's accuracy on MMA is in A.2.

# 4. Experiment

In this section, we conduct extensive experiments to answer the following questions:

1. How well do current leading MLLMs perform on our MMA benchmark, and how significant is the performance difference between MLLMs and human annotation? Sec 4.3.1
2. Explore the reasons why MLLMs lag behind humans in MMA benchmark? Sec 4.3.2
3. How well do the models handle each type of ambiguity? Sec 4.3.4
4. To what extent does model scale (number of parameters) influence performance? Sec 4.3.5

## 4.1. Evaluation models

We evaluate 25 recent multimodal LLMs on our benchmark, including 6 proprietary MLLMs: GPT-4-vision [5], GPT-4o [56], Claude-3-Opus [57], Gemini-1.5-Pro [58], Gemini-1.0-Pro-Vision [59], Claude 3.5 Sonnet [57] and 11 open-source MLLMs: LLaVAV-Next [7], VILA1.5 [60], Yi-VL-34b [61], InternVL-Chat-V1.5[62], InternVL2[62], CogVLM2-LLaMa3-Chat-19B [63], DeepSeek-VL-7b-Chat [64], MiniCPM-Llama3-V-2.5[65], HPT1.5-Air [66], Qwen2-V[67], LLaVA-OneVision[68].

## 4.2. Evaluation metrics

Evaluating the ability of MLLMs to resolve ambiguity in multimodal settings requires metrics that go beyond standard accuracy measures. Due to budgetary constraints, human validation was conducted on a stratified sample of 402 core questions (804 images) from MMA v1. To capture the nuances of this challenge, we introduce this metrics for the MMA benchmark:



Figure 3: **Performance comparison of MLLMs on different ambiguity types**.



Figure 4: **Ablation of parameter size vs. ambiguity accuracy** across ambiguity types.

**Ambiguity Accuracy** ($Amb_A$)  This metric measures the proportion of ambiguous questions for which the model correctly answers for **both** paired images. A high $Amb_A$ score indicates that the model is not defaulting to a single interpretation of the ambiguous text but is instead effectively integrating visual information to provide context-appropriate answers.

Formally, for each ambiguous question $i$ in the benchmark, let $(I_{i,1}, I_{i,2})$ be its paired images, which depict mutually exclusive scenarios. The ground-truth answers for each pair are $(y_{i,1}, y_{i,2})$, and the model's corresponding predictions are $(\hat{y}_{i,1}, \hat{y}_{i,2})$. A score of 1 is assigned to question $i$ only if the model correctly answers for both images:

---

$$\delta_i = \begin{cases} 1, & \text{if } (\hat{y}_{i,1} = y_{i,1}) \wedge (\hat{y}_{i,2} = y_{i,2}) \\ 0, & \text{otherwise.} \end{cases}$$

. The overall Ambiguity Accuracy ($Amb_A$) is the mean of these scores across all $N$ questions in the benchmark: $Amb_A = \frac{1}{N} \sum_{i=1}^{N} \delta_i$ This formulation directly evaluates the model's ability to resolve textual ambiguity using visual context—the core challenge addressed by our benchmark.

## 4.3. Main results

| | Adjective | Noun | Verb | Attachment | Coordination | Structural | Pragmatic | Idiom | Lexical | Syntactic | Semantic | Overall |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Proprietary MLLMs:** | | | | | | | | | | | | |
| Gemini 1.5 Flash [69] | 0.78 | 0.73 | 0.71 | 0.70 | 0.52 | 0.25 | 0.75 | 0.47 | 0.74 | 0.54 | 0.52 | 0.64 |
| Gemini 1.5 Pro [58] | 0.76 | 0.75 | 0.70 | 0.52 | 0.57 | 0.14 | 0.77 | 0.47 | 0.75 | 0.53 | 0.53 | 0.65 |
| Gemini 2.0 Flash [70] | 0.83 | 0.73 | 0.77 | 0.50 | 0.51 | 0.25 | 0.79 | 0.50 | 0.75 | 0.45 | 0.55 | 0.65 |
| Gemini 2.0 Pro [71] | _0.86_ | _0.85_ | _0.80_ | **0.75** | **0.73** | **0.50** | **0.87** | _0.58_ | _0.85_ | **0.70** | 0.63 | _0.75_ |
| GPT-4o-mini [72] | 0.72 | 0.76 | 0.61 | 0.60 | 0.48 | 0.25 | 0.65 | 0.52 | 0.74 | 0.49 | 0.55 | 0.65 |
| GPT-4o-0806[56] | **0.88** | **0.88** | 0.68 | _0.72_ | _0.59_ | 0.38 | _0.81_ | **0.60** | **0.86** | _0.59_ | _0.65_ | **0.76** |
| Claude 3 Opus [57] | 0.73 | 0.56 | 0.38 | 0.00 | 0.16 | 0.00 | 0.25 | 0.16 | 0.57 | 0.08 | 0.21 | 0.38 |
| Claude 3.5 Sonnet [57] | 0.83 | 0.83 | **0.86** | 0.67 | 0.47 | **0.50** | 0.73 | 0.30 | 0.83 | 0.53 | **0.67** | 0.74 |
| Proprietary Average | 0.80 | 0.76 | 0.69 | 0.56 | 0.50 | 0.28 | 0.70 | 0.45 | 0.76 | 0.49 | 0.54 | 0.65 |
| **Open-source MLLMs:** | | | | | | | | | | | | |
| LLaVA-NeXT-34B [7] | 0.87 | 0.80 | 0.5 | 0.08 | 0.59 | 0.00 | 0.40 | 0.41 | 0.79 | 0.33 | 0.40 | 0.60 |
| LLaVA-NeXT-13B [7] | 0.67 | 0.64 | 0.38 | 0.00 | 0.09 | 0.00 | 0.33 | 0.59 | 0.63 | 0.05 | 0.45 | 0.48 |
| LLaVA-NeXT-7B [7] | 0.13 | 0.60 | 0.13 | 0.00 | 0.00 | 0.14 | 0.28 | 0.27 | 0.52 | 0.02 | 0.27 | 0.37 |
| VILA1.5-40b [60] | 0.73 | 0.81 | 0.63 | 0.23 | 0.55 | 0.00 | 0.60 | 0.38 | 0.79 | 0.36 | 0.49 | 0.63 |
| VILA1.5-13b [60] | 0.40 | 0.70 | 0.13 | 0.00 | 0.14 | 0.14 | 0.38 | 0.49 | 0.63 | 0.10 | 0.43 | 0.49 |
| VILA1.5-3b [60] | 0.13 | 0.18 | 0.25 | 0.08 | 0.09 | 0.14 | 0.18 | 0.08 | 0.18 | 0.10 | 0.13 | 0.15 |
| Yi-VL-34b [61] | 0.73 | 0.63 | 0.25 | 0.08 | 0.14 | 0.00 | 0.45 | 0.24 | 0.62 | 0.10 | 0.35 | 0.46 |
| InternVL-Chat-V1-5 [62] | 0.80 | 0.83 | 0.63 | 0.38 | 0.55 | 0.14 | 0.70 | 0.54 | 0.82 | 0.43 | 0.62 | 0.70 |
| InternVL2-40B [62] | 0.60 | 0.60 | 0.50 | 0.15 | 0.59 | _0.43_ | 0.50 | 0.27 | 0.59 | 0.43 | 0.47 | 0.53 |
| Cogvlm2 [63] | 0.33 | 0.57 | 0.13 | 0.00 | 0.36 | 0.00 | 0.38 | 0.43 | 0.52 | 0.19 | 0.40 | 0.43 |
| DeepSeek-VL [64] | 0.47 | 0.70 | 0.50 | 0.23 | 0.27 | 0.00 | 0.53 | 0.38 | 0.66 | 0.21 | 0.45 | 0.53 |
| MiniCPM-Llama3-V 2.5 [65] | 0.00 | 0.12 | 0.25 | 0.15 | 0.14 | 0.00 | 0.23 | 0.05 | 0.11 | 0.12 | 0.14 | 0.12 |
| HPT 1.5 Air [66] | 0.80 | 0.76 | 0.25 | 0.23 | 0.23 | 0.00 | 0.53 | 0.59 | 0.73 | 0.19 | 0.56 | 0.59 |
| Qwen2-VL-72B [67] | 0.79 | 0.72 | 0.50 | 0.40 | 0.41 | **0.50** | 0.58 | 0.10 | 0.72 | 0.42 | 0.51 | 0.61 |
| Qwen2-VL-7B [67] | **0.93** | 0.77 | **0.83** | 0.00 | 0.37 | 0.33 | 0.57 | 0.10 | 0.79 | 0.26 | 0.50 | 0.62 |
| LLaVA-OneVision-72B [68] | **0.93** | 0.61 | 0.50 | **0.77** | 0.59 | 0.14 | 0.41 | 0.00 | 0.63 | 0.57 | 0.35 | 0.54 |
| LLaVA-OneVision-7B [68] | 0.47 | 0.74 | 0.38 | 0.23 | 0.45 | 0.00 | 0.50 | 0.18 | 0.69 | 0.31 | 0.45 | 0.56 |
| Open-sourced Average | 0.58 | 0.63 | 0.39 | 0.18 | 0.33 | 0.12 | 0.44 | 0.30 | 0.61 | 0.25 | 0.41 | 0.50 |
| **Human:** | | | | | | | | | | | | |
| Human Average | 0.83 | 0.93 | 0.83 | 1.00 | 0.90 | 0.63 | 0.82 | 0.98 | 0.91 | 0.89 | 0.85 | 0.89 |

Table 3: **Overall performance comparisons (Amb_A) of MLLMs on different ambiguity types**. The best results are **bold**. The second best results are underlined.

### 4.3.1. Overall performance

As shown in Table 3, the mean ambiguity accuracy (Amb_A) of MLLMs varies significantly across different ambiguity types, highlighting challenges in handling structural and pragmatic ambiguities. However, a clear gap remains when comparing these models to human performance, which significantly outperforms the MLLMs.

Proprietary models, such as Gpt-4o (76%), achieve the best overall performance on Amb_A and excel at handling syntactic ambiguities (59%). Among open-source models, InternVL-Chat-V1-5 (69.7%) shows strong performance, particularly in lexical categories (82%), achieving nearly comparable performance to Gpt-4o.

Despite these advancements, the best-performing models like Gemini 2.0 pro and GPT-4o still show a substantial gap when compared to human performance. Gemini 2.0 pro achieves an overall accuracy of 75%, which is 14% lower than the human benchmark of 89%. Similarly, Gpt-4o performs 13% lower than human performance with an overall accuracy of 76%. Claued 3.5 and InternVL-Chat-V1-5 also underperform humans by 15% and 19%, respectively, with overall accuracy of 74% and 70%. This significant performance gap is particularly evident in tasks involving syntactic and semantic ambiguities. For example, Gemini 2.0 pro and GPT-4o achieve accuracy of 70% and 59% in syntactic ambiguities, respectively, compared to the human accuracy of 89%. In semantic ambiguities, Claude 3.5 Sonnet and Gpt-4o achieve 67% and 65%, respectively, while humans achieve 85%.

### 4.3.2. Explore the reasons for the gap between SOTA models and human

In order to investigate the reasons behind the performance gap between models and humans, we conducted the following experiments:

**MLLMs Performance with Text-Only Input:** Initially, we explored if the inherent **complexity of the tasks or human-crafted questions** might contribute to the performance gap. To this end, we

Table 4: **MLLM Performance with Text-Only Input**. This metric evaluates the ratio of selecting a correct answer when MLLMs receive text-only input. It measures language understanding ability and reduces potential bias from human-crafted questions.

| Model | Attachment | Overall |
|---|---|---|
| Claude 3.5 Sonnet | 0.77 | 0.83 |
| GPT-4 Vision | 1.00 | 0.90 |
| Claude 3 Opus | 1.00 | 0.88 |
| GPT-4o | 0.85 | 0.88 |
| InternVL-Chat-V1.5 | 0.85 | 0.86 |

Table 5: **MLLMs' Error Consistency Rate**. The ratio of cases where MLLMs provide the same answer when presented with two different images, measuring how much models neglect image information in resolving ambiguities.

| Model | Lexical | Syntactic | Semantic | Overall |
|---|---|---|---|---|
| Claude 3 Opus | 0.86 | 0.72 | 0.89 | 0.84 |
| GPT-4o | 0.72 | 0.83 | 0.79 | 0.78 |
| InternVL-1.5 | 0.62 | 0.83 | 0.69 | 0.71 |
| DeepSeek-VL | 0.69 | 0.76 | 0.69 | 0.71 |
| HPT 1.5 Air | 0.66 | 0.82 | 0.74 | 0.74 |
| VILA1.5-40b | 0.73 | 0.78 | 0.95 | 0.83 |
| Yi-VL-34b | 0.65 | 0.84 | 0.86 | 0.77 |

assessed the accuracy of MLLMs when they were provided solely with text inputs. The metric used represents the rate at which the model's response matches one of the correct answers in each pair of data (ambiguity pair), it is considered accurate. As shown in Table 4, MLLMs demonstrate high accuracy when provided with only text input. The overall accuracy rates range from 83% to 90%, with GPT-4 Vision achieving the highest at 90%. Notably, performance is consistently strong across lexical, syntactic, and semantic categories, with most models scoring above 80% in each. Claude 3.5 Sonnet shows the most balanced performance across categories, while others like InternVL-Chat-V1-5 exhibit some variability (e.g., 90% lexical vs. 74% syntactic). These results indicate that minor textual issues have minimal impact on MLLMs' ability to select correct answers.

**MLLMs' Error Consistency Rate:** This Error Consistency Rate (ECR) - defined as the rate of selecting the same answer among incorrect cases. As shown in Table 5, when MLLMs made errors, they demonstrated a high consistency rate in choosing the same option twice. This rate ranged from 71% to 84% overall, depending on the model. The consistently high rates across lexical, syntactic, and semantic levels indicate that these models often failed to effectively leverage visual information when answering questions. Instead, they exhibited a strong bias towards the text modality, relying primarily on textual cues even when visual information was available. More error analysis are in Appendix F.1.

In summary, the experimental results clearly indicate that the performance gap between MLLMs and humans does not stem from the inherent complexity of the tasks or the construction of the questions, as evidenced by the high accuracy rates with text-only inputs. Rather, the persistent performance gap is largely due to the models' failure to adequately process and integrate visual information to clarify the textual ambiguity. The tendency of MLLMs to repeat the same answers, even when presented with different visual contexts, highlights a pronounced bias towards textual information instead of leveraging visual information.

### 4.3.3. Gap between proprietary models and open-sourced models

On average, proprietary models demonstrate better performance than open-sourced models in the MMA task. Specifically, proprietary models achieve 65% in Amb_A, while open-sourced models obtain 49.47% in Amb_A as Table 3 shows. For both indicators, proprietary models outperform open-sourced models.

### 4.3.4. Syntactic ambiguity and semantic ambiguity are more challenging than lexical ambiguity

For both Amb_A, all models (except MiniCPM) perform better on lexical ambiguity and worse on syntactic and semantic ambiguities (Figure 5). Lexical ambiguity, which involves straightforward word meanings, is easier for models to handle. For example, InternVL-Chat-V1-5 achieves an accuracy of 82% on lexical ambiguities, significantly higher than its performance on syntactic (43%) and semantic (62%) ambiguities. This trend is consistent across most models; for instance, Gpt-4o shows 86% accuracy on lexical ambiguities but drops to 59% and 65% on syntactic and semantic ambiguities, respectively.

Figure 5: All models except MiniCPM-Llama3-V 2.5 perform better on Lexical ambiguity than Syntactic ambiguity and semantic ambiguity for Amb_A.

Syntactic ambiguities present a unique challenge because they involve the relationships between components within a sentence. Often, even a short modifier can introduce ambiguity, making it difficult for models to resolve these cases without fine-grained analysis. To effectively handle syntactic ambiguities, models need not only a more granular approach to language processing but also the capability to accurately recognize positional relationships and details in images. This requires a higher level of precision compared to lexical and semantic ambiguities. Models like Cogvlm2 and VILA-3b, for example, perform poorly in this category, with accuracies of 19% and 10%, respectively. Similarly, semantic ambiguities, which involve nuanced meanings and context, are also difficult for models to resolve. For instance, VILA-40b achieves only 49% accuracy on semantic ambiguities, despite a higher performance on lexical (79%).

### 4.3.5. Scaling experiment on MMA

Our investigation into the impact of parameter count on the MMA benchmark reveals inconsistent performance trends across different models. While some model series, such as VILA1.5, demonstrate a clear improvement in ambiguity accuracy with increased parameter counts Figure 4, this scaling benefit is not universal. Notably, we observed instances where larger models performed worse than their smaller counterparts (like Qwen2-VL-72B vs. 7B, LLaVA-OneVision-72B vs. 7B). This suggests that resolving visual ambiguity does not uniformly improve with more parameters; larger models might overfit to common visual-language associations from training, potentially diminishing their ability to handle nuanced ambiguities compared to smaller models that might retain better uncertainty representation for this specific task. Furthermore, Performance regressions in newer variants (like InternVL) show that our benchmark measures abilities not automatically improved through standard model iteration or data scaling. This indicates that parameter scaling alone is insufficient—dataset design, data collection, and architectural or training strategies play crucial roles in achieving strong performance. MMA thus helps identify specialized skills requiring targeted development beyond general scaling..

## 5. Conclusion

This paper introduces MMA, the first benchmark designed specifically to evaluate the ability of MLLMs to understand and respond to ambiguous queries. MMA leverages a multiple-choice visual question-answering format, presenting MLLMs with a question and two images depicting contrasting scenarios that lead to different correct answers. Our evaluation of 25 MLLMs, including both limited-access and open-sourced models, reveals a significant performance gap compared to human performance. While humans achieve an accuracy of 88.97%, the MLLMs average only 54.8% accuracy. This indicates a fundamental challenge for current MLLMs: effectively integrating scenario-specific visual information to disambiguate questions and arrive at the correct answer. Even the top-performing model, Gpt-4o, attains only about 76% accuracy, highlighting considerable room for developing MLLMs that can effectively leverage visual information to clarify the textual ambiguity and capable of human-level understanding and reasoning in complex, real-world scenarios.

# References

[1] Matthew Turk. Multimodal interaction: A review. *Pattern recognition letters*, 36:189–195, 2014.

[2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[3] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the conference. Association for computational linguistics. Meeting*, volume 2019, page 6558. NIH Public Access, 2019.

[4] Peng Xu, Xiatian Zhu, and David A. Clifton. Multimodal learning with transformers: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10):12113–12132, 2023. doi: 10.1109/TPAMI.2023.3275156.

[5] OpenAI. Gpt-4 technical report, 2024.

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

[7] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL `https://llava-vl.github.io/blog/2024-01-30-llava-next/`.

[8] Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. Mobile-agent: Autonomous multi-modal mobile device agent with visual perception. *arXiv preprint arXiv:2401.16158*, 2024.

[9] Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei Lin, Saravan Rajmohan, et al. Ufo: A ui-focused agent for windows os interaction. *arXiv preprint arXiv:2402.07939*, 2024.

[10] Keen You, Haotian Zhang, Eldon Schoop, Floris Weers, Amanda Swearngin, Jeffrey Nichols, Yinfei Yang, and Zhe Gan. Ferret-ui: Grounded mobile ui understanding with multimodal llms. *arXiv preprint arXiv:2404.05719*, 2024.

[11] Yutao Cheng, Zhao Zhang, Maoke Yang, Hui Nie, Chunyuan Li, Xinglong Wu, and Jie Shao. Graphic design with large multimodal model. *arXiv preprint arXiv:2404.14368*, 2024.

[12] Jieru Lin, Danqing Huang, Tiejun Zhao, Dechen Zhan, and Chin-Yew Lin. Design-probe: A graphic design benchmark for multimodal large language models. *arXiv preprint arXiv:2404.14801*, 2024.

[13] Jie Gao, Simret Araya Gebreegziabher, Kenny Tsu Wei Choo, Toby Jia-Jun Li, Simon Tangi Perrault, and Thomas W Malone. A taxonomy for human-llm interaction modes: An initial exploration. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, pages 1–11, 2024.

[14] Soufia Bahmani. Enhancing human-computer interaction through large language models: Opportunities, challenges, and future directions.

[15] Yi Yang, Qingwen Zhang, Ci Li, Daniel Simões Marta, Nazre Batool, and John Folkesson. Human-centric autonomous systems with llms for user command reasoning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 988–994, 2024.

[16] Thomas Wasow, Amy Perfors, and David Beaver. The puzzle of ambiguity. *Morphology and the web of grammar: Essays in memory of Steven G. Lapointe*, pages 265–282, 2005.

[17] Alisa Liu, Zhaofeng Wu, Julian Michael, Alane Suhr, Peter West, Alexander Koller, Swabha Swayamdipta, Noah Smith, and Yejin Choi. We're afraid language models aren't modeling ambiguity. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 790–807, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.51. URL `https://aclanthology.org/2023.emnlp-main.51/`.

[18] Hyuhng Joon Kim, Youna Kim, Cheonbok Park, Junyeob Kim, Choonghyun Park, Kang Min Yoo, Sang goo Lee, and Taeuk Kim. Aligning language models to explicitly handle ambiguity, 2024.

[19] Sigrid Norris. *Analyzing multimodal interaction: A methodological framework.* Routledge, 2004.

[20] Mohammad Taher Pilehvar and Jose Camacho-Collados. WiC: the word-in-context dataset for evaluating context-sensitive meaning representations. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1267–1273, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1128. URL `https://aclanthology.org/N19-1128/`.

[21] Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. Conll-2012 shared task: Modeling multilingual unrestricted coreference in ontonotes. In *Joint conference on EMNLP and CoNLL-shared task*, pages 1–40, 2012.

[22] Davide Buscaldi, Anne-Kathrin Schumann, Behrang Qasemizadeh, Haïfa Zargayouna, and Thierry Charnois. Semeval-2018 task 7: Semantic relation extraction and classification in scientific papers. In *International Workshop on Semantic Evaluation (SemEval-2018)*, pages 679–688, 2017.

[23] Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. AmbigQA: Answering ambiguous open-domain questions. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5783–5797, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.466. URL `https://aclanthology.org/2020.emnlp-main.466/`.

[24] Jonathan Pilault, Xavier Garcia, Arthur Bražinskas, and Orhan Firat. Interactive-chain-prompting: Ambiguity resolution for crosslingual conditional generation with interaction. In Jong C. Park, Yuki Arase, Baotian Hu, Wei Lu, Derry Wijaya, Ayu Purwarianti, and Adila Alfa Krisnadhi, editors, *Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 455–483, Nusa Dua, Bali, November 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.ijcnlp-main.31. URL `https://aclanthology.org/2023.ijcnlp-main.31/`.

[25] Yuewei Yuan, Chaitanya Malaviya, and Mark Yatskar. AmbiCoref: Evaluating human and model sensitivity to ambiguous coreference. In Andreas Vlachos and Isabelle Augenstein, editors, *Findings of the Association for Computational Linguistics: EACL 2023*, pages 1023–1030, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-eacl.75. URL `https://aclanthology.org/2023.findings-eacl.75/`.

[26] Yevgeni Berzak, Andrei Barbu, Daniel Harari, Boris Katz, and Shimon Ullman. Do you see what I mean? visual resolution of linguistic ambiguities. In Lluís Màrquez, Chris Callison-Burch, and Jian Su, editors, *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1477–1487, Lisbon, Portugal, September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1172. URL `https://aclanthology.org/D15-1172/`.

[27] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, and Feng Zhao. Are we on the right way for evaluating large vision-language models?, 2024.

[28] Jiwan Chung, Seungwon Lim, Jaehyun Jeon, Seungbeen Lee, and Youngjae Yu. Can visual language models resolve textual ambiguity with visual cues? let visual puns tell you!, 2024. URL https://arxiv.org/abs/2410.01023.

[29] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, Dinesh Manocha, and Tianyi Zhou. Hallusionbench: An advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models, 2024. URL https://arxiv.org/abs/2310.14566.

[30] Fuwen Luo, Chi Chen, Zihao Wan, Zhaolu Kang, Qidong Yan, Yingjie Li, Xiaolong Wang, Siyu Wang, Ziyue Wang, Xiaoyue Mi, Peng Li, Ning Ma, Maosong Sun, and Yang Liu. Codis: Benchmarking context-dependent visual comprehension for multimodal large language models, 2024. URL https://arxiv.org/abs/2402.13607.

[31] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi, 2023.

[32] Xiaolong Wang, Zhaolu Kang, Wangyuxuan Zhai, Xinyue Lou, Yunghwei Lai, Ziyue Wang, Yawen Wang, Kaiyu Huang, Yile Wang, Peng Li, and Yang Liu. MUCAR: Benchmarking multilingual cross-modal ambiguity resolution for multimodal large language models. In Christos Christodoulopoulos, Tanmoy Chakraborty, Carolyn Rose, and Violet Peng, editors, *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 15026–15048, Suzhou, China, November 2025. Association for Computational Linguistics. ISBN 979-8-89176-332-6. doi: 10.18653/v1/2025.emnlp-main.760. URL https://aclanthology.org/2025.emnlp-main.760/.

[33] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*, 2019.

[34] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning, 2020.

[35] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32, 2019.

[36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

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

[38] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023.

[39] Xinyu Ma, Xuebo Liu, Derek F. Wong, Jun Rao, Bei Li, Liang Ding, Lidia S. Chao, Dacheng Tao, and Min Zhang. 3AM: An ambiguity-aware multi-modal machine translation dataset. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, Torino, Italia, May 2024. ELRA and ICCL. URL `https://aclanthology.org/2024.lrec-main.1`.

[40] Yihang Li, Shuichiro Shimizu, Weiqi Gu, Chenhui Chu, and Sadao Kurohashi. Visa: An ambiguous subtitles dataset for visual scene-aware machine translation. *arXiv preprint arXiv:2201.08054*, 2022.

[41] Tinghui Zhu, Qin Liu, Fei Wang, Zhengzhong Tu, and Muhao Chen. Unraveling cross-modality knowledge conflicts in large vision-language models. *arXiv preprint arXiv:2410.03659*, 2024.

[42] 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.

[43] Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7w: Grounded question answering in images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4995–5004, 2016.

[44] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017.

[45] 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.

[46] 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.

[47] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.

[48] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

[49] Elias Stengel-Eskin, Jimena Guallar-Blasco, Yi Zhou, and Benjamin Van Durme. Why did the chicken cross the road? rephrasing and analyzing ambiguous questions in vqa. *arXiv preprint arXiv:2211.07516*, 2022.

[50] Stephen Ullmann. *Semantics: An Introduction to the Science of Meaning*. Basil Blackwell / Oxford : Blackwell, 1962. First edition.

[51] John Lyons. *Semantics (Vol. 1 & 2)*. Cambridge University Press, 1977. doi: 10.1017/CBO9780511620614.

[52] James Allen. *Natural Language Understanding*. Benjamin/Cummings series in computer science. Benjamin/Cummings, 1995. ISBN 0805303340.

[53] Daniel Jurafsky and James H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Prentice Hall, Upper Saddle River, NJ, 1st edition, 2000. ISBN 9780130950697.

[54] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.

[55] OpenAI. Dall-e-3, 2024. URL `https://openai.com/index/dall-e-3/`.

[56] OpenAI. Hello gpt-4o, 2024. URL `https://openai.com/index/hello-gpt-4o/`.

[57] Anthropic. Claude 3: A new generation of ai, 2024. URL `https://docs.anthropic.com/en/docs/models-overview#claude-3-a-new-generation-of-ai`.

[58] Google DeeoMind. Gemini pro1.5, 2024. URL `https://deepmind.google/technologies/gemini/pro/`.

[59] Google DeeoMind. Gemini pro vision, 2023. URL `https://deepmind.google/technologies/gemini/pro/`.

[60] Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models, 2023.

[61] 01. AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. Yi: Open foundation models by 01.ai, 2024.

[62] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, Ji Ma, Jiaqi Wang, Xiaoyi Dong, Hang Yan, Hewei Guo, Conghui He, Botian Shi, Zhenjiang Jin, Chao Xu, Bin Wang, Xingjian Wei, Wei Li, Wenjian Zhang, Bo Zhang, Pinlong Cai, Licheng Wen, Xiangchao Yan, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhai Wang. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites, 2024.

[63] Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvlm: Visual expert for pretrained language models, 2023.

[64] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. Deepseek-vl: Towards real-world vision-language understanding, 2024.

[65] OpenBMB. Minicpm-llama3-v 2.5, 2024. URL `https://github.com/OpenBMB/MiniCPM-V?tab=readme-ov-file#minicpm-llama3-v-25`.

[66] HYPERGAI. Hpt 1.5 air: Best open-sourced 8b multimodal llm with llama 3, 2024. URL `https://www.hypergai.com/blog/hpt-1-5-air-best-open-sourced-8b-multimodal-llm-with-llama-3`.

[67] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.

[68] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.

[69] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.

[70] DeepMind. Gemini flash, n.d.. URL `https://deepmind.google/technologies/gemini/flash/`. Accessed: 2025-02-16.

[71] DeepMind. Gemini pro, n.d.. URL `https://deepmind.google/technologies/gemini/pro/`. Accessed: 2025-02-16.

[72] OpenAI. gpt-4o-mini-advancing-cost-efficient-intelligence, 2024. URL `https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/`.

# Appendix

## A. Limitation

**Data collection** Due to constraints on the number of participants, the dataset size is limited in certain categories. Despite this limitation, we emphasize that the quality and representativeness of the dataset are more crucial for establishing a meaningful benchmark than merely the number of samples. As demonstrated in Table 3, the considerable performance discrepancy between human participants and MLLM responses underscores the benchmark's effectiveness in highlighting the current challenges that MLLMs face, particularly their inability to adequately utilize visual context to resolve textual ambiguities. Moving forward, we are committed to expanding the dataset in future iterations of the benchmark, aiming to broaden its scope and enhance its representational validity.

**Question design** In our benchmark, both images and texts are designed to provide context information to model the multi-modal real-world cases. Due to the paper presentation problem, how to present some questions naturally presents certain challenges. We conducted experiments with text-only input and found that MLLMs demonstrate high accuracy, ranging from 83% to 90% (as Table 4 shows). However, when errors occurred, models consistently chose the same incorrect answers (as Table 5 shows). These results clearly indicate MLLMs have a strong bias towards text-based information and a failure to effectively incorporate visual context.

**Real-world likeness** Some of images used in our benchmark are generated by generative models. The images in our benchmark are specifically chosen to provide the necessary context to clarify ambiguities in the accompanying texts. Due to the current limitations of search engines, which struggle with semantic search, it is challenging to find suitable images that naturally align with the required context (This doesn't mean that these images don't exist.). Therefore, using generated images is the most effective approach. They are instrumental in simulating the diverse and often unconventional situations that MLLMs encounter in real applications. MLLMs are expected to perform comparably to humans in these scenarios, regardless of the variability in inputs. However, our human study shows that humans can achieve approximately 90% accuracy on this benchmark without any additional interactions. This sharply contrasts with the average accuracies of 58% for closed-source models and 47% for open-source models.

## B. Distribution of Dataset

As shown in 3, the MMA dataset consists of 804 images and 402 questions, covering three main types of ambiguity: lexical ambiguity, syntactic ambiguity, and semantic ambiguity. These main categories are further divided into eight sub-categories: noun ambiguity, verb ambiguity, and adjective ambiguity (under lexical ambiguity); attachment ambiguity, coordination ambiguity, and structural ambiguity (under syntactic ambiguity); and pragmatic ambiguity and idiomatic ambiguity (under semantic ambiguity). For the ablation study, we designed 120 novel question-answer pairs based on the same images from the Noun category.

Regarding data proportions, while lexical ambiguities are significantly more common in everyday language, syntactic and semantic ambiguities, though equally important, are less frequent. Syntactic ambiguities, often arising from complex grammatical structures, are inherently less common.

Therefore, the distribution of lexical, syntactic, and semantic ambiguities in our dataset is designed to reflect their natural prevalence in real-world communication. Our dataset's proportions aim to mirror this natural distribution, ensuring the benchmark's realism and applicability to real-world scenarios.



Figure 6: Ambiguity Type Composition of MMA benchmark

## C. Benchmark and Evaluation Resources

To facilitate benchmarking, we've made the dataset available.

For evaluation purposes, you can utilize the code provided in our github webpage.

## D. Image Usage and Copyright Claims

Our images are sourced from search engines (such as Google and Bing) and text-to-image models (such as Stable-Diffusion and DALL-E). All collected images are used exclusively to support our non-profit research project, MMA Benchmark. If you own the copyright to any images used in this project and believe that your rights have been violated, please contact us. We are willing to compensate for the usage of your images.

## E. Future Work

**Additional Modalities** The world is multimodal rather than just bimodal. For instance, audio plays an important role in daily life, and there are some ambiguities caused by audio. For example, the phrases "He's a great **rapper**" and "He's a great **wrapper**" sound similar but refer to completely different things. With a concrete scene provided, the meaning of a segment of audio can be uniquely determined.

**Additional Languages** Language-specific features and rhetorical devices vary widely, influencing how information is processed and understood. For instance, the use of 'Huwen' in ancient Chinese literature requires an understanding of how meanings are intricately split and reconnected across sentences. Expanding MLLMs to accommodate the linguistic structures and subtleties of various languages could improve their applicability and accuracy in global communication contexts. This development would necessitate models that are not only multilingual but also sensitive to cultural and contextual nuances within languages.

**Multiple Images per Sentence** Lexical ambiguities can extend beyond dual interpretations, with some words or phrases having multiple meanings. Current models often limit context to one or two visual representations per sentence. By providing multiple images that correspond to each

potential meaning of a sentence, MLLMs can be trained to discern finer distinctions in word usage and context. This enhancement would allow models to handle more complex scenarios where multiple interpretations are valid, reflecting the true complexity of human language and cognition.

# F. Ablation study

**Same images with lexical or semantic questions**  To understand why MLLMs perform better on lexical ambiguity compared to semantic ambiguity, we explored how changing the question type on noun ambiguity impacts their performance. We created two versions of questions for noun categories: the first being the most direct, "What's the meaning of <Noun>?", and the second incorporating reasoning into the question. For example, given an image of a table, a synonym question for lexical ambiguity might be "What is the meaning of table?" where the model identifies "table" as a piece of furniture. In contrast, a reasoning question for semantic ambiguity would be "How can we best utilize the space on this table?" which requires the model to consider various uses of the table. This type of question tests the model's ability to perform object grounding and higher-order reasoning, areas where MLLMs often show weaker performance due to their reliance on pattern recognition rather than true comprehension. More examples are given in Appendix.

As Figure 7 shows, GPT-4 Vision performs well on noun word ambiguity with a score of 90% but drops to 59% on noun reasoning ambiguity. Similarly, Gemini-1.5 shows a significant drop from 83% in noun word ambiguity to 63% in noun reasoning ambiguity. Intern-VL-Chat-V1-5, while achieving 92% in noun word ambiguity, sees a decline to 75% in noun reasoning ambiguity. These examples highlight the challenges MLLMs face in understanding and reasoning about more complex and context-dependent scenarios.



Figure 7: **The performance comparison for question types**, where The Noun_word refers to questions that solely inquire about the meaning of a noun word, while the Noun_reasoning involves questions that require the reasoning ability to answer. The details and examples are given in Appendix.

**Prompt Ablation Study**  To examine the influence of prompt phrasing on model behavior, we conducted an ablation study comparing our original concise evaluation prompt, "Select the most suitable option for the question from the following choices," with a more explicit variant, "Choose the description that most matches the image." The experiment was performed on the Qwen-VL-7B model to assess whether explicit instructions improve its ability to utilize visual context for disambiguation, as summarized in Table 6. The results show no consistent gain for the explicit prompt: while categories such as Coordination, Structural, Pragmatic, Idiom, and Syntactic ambiguity exhibit noticeable gains, others—including Adjective, Verb, and Lexical ambiguity—decline substantially. Consequently, the overall performance remains unchanged (0.62 → 0.62), indicating that Qwen-VL-7B does not reliably benefit from more explicit prompting when resolving ambiguity.

## F.1. Error Analysis

Errors can be categorized into three main types: **uni-modal image issues, uni-modal text issues, and cross-modal text bias.** An analysis of the error distribution in GPT-4o reveals that cross-modal text bias errors constitute the majority of all errors(see Figure 8). This finding suggests that there is significant room for improvement MMA benchmark.

Table 6: **Prompt ablation results.** Comparison between the original concise and explicit prompts on Qwen-VL-7B.

| Category | Original | Explicit |
|---|---|---|
| Adjective | 0.93 | 0.67 |
| Noun | 0.77 | 0.75 |
| Verb | 0.83 | 0.63 |
| Attachment | 0.00 | 0.00 |
| Coordination | 0.37 | 0.55 |
| Structural | 0.33 | 0.43 |
| Pragmatic | 0.57 | 0.60 |
| Idiom | 0.10 | 0.49 |
| Lexical | 0.79 | 0.74 |
| Syntactic | 0.26 | 0.36 |
| Semantic | 0.50 | 0.54 |
| **Overall (Amb_A)** | **0.62** | **0.62** |



Figure 8: **Error type distribution of GPT-4o**, where we see cross-model text bias accounts for half of the cases.



Figure 9: The ablation study about the parameter number and the ambiguity accuracy performance on different ambiguity types.

**Uni-modal Image Issues (22.1%)** In this type of error, the model fails to capture the essential information conveyed by the image. To address this issue, visual prompts, such as red bounding boxes, can be incorporated to redistribute the attention of the Multimodal Large Language Model (MLLM). By emphasizing the crucial elements of the image, the model can be guided towards generating the correct answer based on the key visual information(see Figure 10).

**Uni-modal Text Issues (27.9%)** In this type of error, the model successfully captures the essential information from the image but provides an incorrect answer due to misinterpreting the text options. To resolve this issue, text prompts can be introduced to guide the MLLMs towards a proper understanding of the textual content. By ensuring accurate comprehension of the text, these prompts can help the model arrive at the correct answer (see Figure 11).

18

Figure 10: Uni-modal Image Issues: the model fails to capture the essential information conveyed by the image.



Figure 11: Uni-modal Text Issues: the model successfully captures the essential information from the image but provides an incorrect answer due to misinterpreting the text options.

**Cross-modal Text Bias (50.0%)** In this category of errors, the model successfully identifies the essential information in the image and comprehends the text options. However, it provides an incorrect answer due to overlooking certain aspects of the visual information while overemphasizing the textual information. To mitigate this issue, text prompts can be introduced to rebalance the attention between the image and text modalities(see Figure 12). By adjusting the relative importance of visual and textual cues, the model can be encouraged to arrive at the correct answer by considering all relevant information from both modalities.

# G. Human Evaluation

To validate our dataset and assess the performance difference between humans and models, we invited five people to participate in benchmark testing. As shown in the table, for each sub-ambiguity class, at least one person achieves an ambiguity accuracy of over 90%, with the exception of Pragmatic ambiguity, where the highest accuracy is 88%. These results demonstrate that our dataset is well-constructed and solvable by humans, serving as a strong validation of the dataset's quality and the feasibility of the task. Humans may fail to answer questions correctly due to a lack of knowledge (such as not understanding the meaning of an idiom), being confused by misleading or similar

Figure 12: Cross-modal Text Bias: the model successfully captures the essential information from the image but provides an incorrect answer due to misinterpreting the text options.

| | Adjective (30) | Noun (238) | Verb (16) | Attachment (24) | Coordination (46) | Structural (14) | Pragmatic (132) | Idiom (22) | Lexical (284) | Syntactic (84) | Semantic (154) | Overall (522) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Person1 | 0.60 | 0.88 | 0.88 | 1.00 | 0.77 | 0.00 | 0.74 | 0.91 | 0.85 | 0.71 | 0.77 | 0.80 |
| Person2 | 0.93 | 0.97 | 1.00 | 1.00 | 0.86 | 1.00 | 0.83 | 1.00 | 0.96 | 0.93 | 0.86 | 0.93 |
| Person3 | 0.80 | 0.94 | 0.50 | 1.00 | 0.91 | 0.71 | 0.88 | 1.00 | 0.90 | 0.90 | 0.90 | 0.90 |
| Person4 | 0.93 | 0.93 | 1.00 | 1.00 | 0.95 | 0.71 | 0.85 | 1.00 | 0.94 | 0.93 | 0.87 | 0.92 |
| Person5 | 0.87 | 0.93 | 0.75 | 1.00 | 1.00 | 0.71 | 0.82 | 1.00 | 0.92 | 0.95 | 0.84 | 0.90 |

Table 7: Five people have different performance across different types of ambiguities

answer options with subtle differences, or struggling to correlate images with text (particularly when the text contains advanced expressions or extended meanings). Structural" type of questions often involve verbs that can function both transitively and intransitively, each usage yielding a distinct meaning. As a result, these questions require more complex grammatical reasoning and the ability to discern subtle syntactic cues. Even native speakers may struggle with certain advanced constructions.

# H. Case Study

More examples of GPT-4o will be presented in this section.

**Example of GPT-4o**

1. Example of Coordination Ambiguity
2. Example of Adjective Ambiguity
3. Example of Noun Ambiguity
4. Example of Structural Ambiguity
5. Example of Attachment Ambiguity
6. Example of Idiomatic Ambiguity
7. Example of Pragmatic Ambiguity
8. Example of Verb Ambiguity

# I. Ethics Statement

This work was conducted in accordance with ethical standards for data collection, human partici-pation, and responsible AI research. All human annotators participated voluntarily with informed

**Question:** The children are building a sandcastle. From the image I offered to you, What can you infer about the children? <Image1>

**Option:**
A.Not specific.
B.These children are building the different sandcastle respectively
C.These children are building the same sandcastle.
D.These children are not doing anything relevant to building sandcastle

<Image1>

**GPT-4o:**
C | Based on the image, the children appear to be working together to build a single sandcastle on the beach. They are both interacting with the same sand structure, suggesting they are collaborating to construct one shared sandcastle rather than building separate ones.

**Ground Truth:**
C.These children are building the same sandcastle.

**Question:** The children are building a sandcastle. From the image I offered to you, What can you infer about the children? <Image2>

**Option:**
A.Not specific.
B.These children are building the different sandcastle respectively
C.These children are building the same sandcastle.
D.These children are not doing anything relevant to building sandcastle

<Image2>

**GPT-4o:**
C: These children are building the same sandcastle
The children in the images are engaged in building a sandcastle, which fits the context of them working on the same sandcastle together.

**Ground Truth:**
B.These children are building the different sandcastle respectively

Figure 13: Example of Coordination Ambiguity



**Question:** "It is sharp." From the image I offered to you, Can you explain the meaning "sharp" here? <Image1>

**Option:**
A.Distinct in outline or detail; clearly defined.
B.Having or showing speed of perception, comprehension, or response.
C.Having an edge or point that is able to cut or pierce something.
D.Sudden and marked.

<Image1>

**GPT-4o:**
C | The image shows a close-up view of hands using a sharp knife to cut into what appears to be a piece of leather or a thick, rigid material. In this context, "sharp" most likely refers to the knife having an edge or point that is able to easily cut or pierce through the material being worked on.

**Ground Truth:**
C.Having an edge or point that is able to cut or pierce something.

**Question:** "It is sharp." From the image I offered to you, Can you explain the meaning "sharp" here? <Image2>

**Option:**
A.Distinct in outline or detail; clearly defined.
B.Having or showing speed of perception, comprehension, or response.
C.Having an edge or point that is able to cut or pierce something.
D.Sudden and marked.

<Image2>

**GPT-4o:**
C: Having an edge or point that is able to cut or pierce something. | explanation: In the context of the image, which appears to depict a stock market chart with a sharp drop in price, "sharp" most suitably refers to a steep or dramatic decline. This usage aligns with the idea of something being "sharp" in the sense of having a sudden, marked change, similar to how a sharp object has a distinct edge or point.

**Ground Truth:**
D.Sudden and marked.

Figure 14: Example of Adjective Ambiguity



**Question:** Can the table be used to place object?<Image1>

**Option:**
A.Yes, we can physical place objects on the wooden table.
B.No, we can't place data structure objects on the wooden table.
C.Yes, we can place data structure objects in the database table.
D.No, we can place physical place objects in the database table.

<Image1>

**GPT-4o:**
A | Yes, we can physically place objects on the wooden table. The image shows a sturdy round wooden table with a flat surface, which would be suitable for placing and holding various physical objects on top of it.

**Ground Truth:**
A.Yes, we can physical place objects on the wooden table.

**Question:** Can the table be used to place object?<Image2>

**Option:**
A.Yes, we can physical place objects on the wooden table.
B.No, we can't place data structure objects on the wooden table.
C.Yes, we can place data structure objects in the database table.
D.No, we can place physical place objects in the database table.

<Image2>

**GPT-4o:**
B: No, we can't place data structure objects on the wooden table.
Explanation: The image provided depicts database tables, which are structures used to store and manage data in a database. These tables cannot be used to physically place objects.

**Ground Truth:**
C.Yes, we can place data structure objects in the database table.

Figure 15: Example of Noun Ambiguity

consent, and no personally identifiable or sensitive information was collected. The dataset was

**Figure 16: Example of Structural Ambiguity**



**Figure 17: Example of Attachment Ambiguity**



**Figure 18: Example of Coordination Ambiguity**

carefully reviewed to exclude harmful, offensive, or privacy-related content. All images were used

Figure 19: Example of Idiomatic Ambiguity



Figure 20: Example of Verb Ambiguity

strictly for non-commercial, academic purposes, as further detailed in Appendix D. The MMA benchmark is released to promote transparent, reproducible, and fair evaluation of multimodal large language models under ambiguous contexts. We emphasize the responsible use of this resource for advancing research in safe and interpretable AI systems.

# J. Usage of Large Language Models

**Data Generation:** Generative AI tools (Like Dall-E, Stable diffusion and other image generation tools) were used to synthesize a few subset of the images in the MMA benchmark to construct controlled ambiguity scenarios.

**Writing Assistance:** Large Language Models were used for language polishing, grammar refinement, and improving readability of the manuscript.

LLMs were not involved in the scientific ideation, interpretation of results, or the formulation of claims. The authors verified all data and take full responsibility for the content.