000 MMA: BENCHMARKING MULTI-MODAL LARGE LAN-001 GUAGE MODELS IN AMBIGUITY CONTEXTS 002 003

Anonymous authors

004

010

Paper under double-blind review

ABSTRACT

011 Multi-Modal Large Language Models (MLLMs) recently demonstrated strong capabilities in both instruction comprehension and responding, positioning them as 012 promising tools for human-computer interaction. However, the inherent ambiguity 013 of language poses a challenge, potentially leading models astray in task implemen-014 tation due to differing interpretations of the same text within varying contexts. In 015 multi-modal settings, visual information serves as a natural aid in disambiguating 016 such scenarios. In this paper, we introduce the first benchmark specifically de-017 signed to evaluate the performance of MLLMs in Ambiguous contexts (MMA). 018 This benchmark employs a multiple-choice visual question-answering format and 019 includes 261 textual contexts and questions with ambiguous meaning. Each question is linked to a pair of images that suggest divergent scenarios, thus leading 021 to different answers given the same question. These questions are stratified into three categories of ambiguity: lexical, syntactic, and semantic, to facilitate a detailed examination of MLLM performance across varying levels of ambiguity. By 023 evaluating 24 proprietary and open-sourced MLLMs, we find that: (1) MLLMs often overlook scenario-specific information provided by images to clarify the 025 ambiguity of texts. When presented with two different contextual images and 026 asked the same question, MLLMs achieved an accuracy rate of only 53.22% in answering both correctly, compared to human performance at 88.97%. (2) Among 028 the three types of ambiguity, models perform best under lexical ambiguity and 029 worst under syntactic ambiguity. (3) Open-sourced models generally perform significantly lower than proprietary MLLMs, with an average performance gap of 031 12.59%, Claude 3.5 Sonnet, emerges as the top model, achieving 74.32% accuracy. 032 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. 034

- INTRODUCTION 1
- 036

Our interaction with the world is inherently multimodal, involving the reception and processing of information across modalities (Turk, 2014). By training on large-scale datasets, multimodal large language models (MLLMs) built-up on transformers (Vaswani et al., 2017; Tsai et al., 2019; Xu 040 et al., 2023), such as GPT-4V (OpenAI, 2024b), Gemini (Team et al., 2023) and LLaVA (Liu et al., 041 2024), have demonstrated strong understanding, reasoning, and even coding ability across vision and 042 language modalities. With visual and language understanding abilities, the realization of MLLM-043 based agents has become feasible, sparking the potential for a variety of innovative applications, such 044 as mobile-operation (Wang et al., 2024a; Zhang et al., 2024; You et al., 2024) and graphics design (Cheng et al., 2024; Lin et al., 2024). These applications highlight the transformative potential of MLLMs in future human-computer interaction (Gao et al., 2024; Bahmani; Yang et al., 2024). 046

047 However, clarity during interactions is not always guaranteed. Ambiguity, which refers to cases where 048 an expression conveys multiple denotations (Wasow et al., 2005; Liu et al., 2023b; Kim et al., 2024), is inherently present in human interactions (Norris, 2004). 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 051 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. 052 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



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.

072 texts. If the model cannot handle ambiguity effectively, there is a risk of misinterpreting the 073 user's original intent, potentially harming the model's reliability. In multimodal contexts, while 074 visual cues provide additional layers of meaning, the capability of MLLMs to effectively manage 075 such ambiguity remains untested. This introduces significant concerns regarding the robustness and 076 reliability of MLLMs, which are essential for their practical deployment.

077 To systematically evaluate and enhance MLLM capabilities in handling these challenges, we introduce a novel benchmark, MLLMs with Ambiguous questions (MMA). This benchmark is specifically 079 designed in a multiple-choice visual-question answering format, featuring 261 questions that each link to a pair of images depicting divergent scenarios. This design ensures that the same question 081 may elicit different correct responses depending on the provided contextual image, thereby testing 082 the model's ability to navigate ambiguity effectively.

083 The questions within MMA are categorized into three types of ambiguity—lexical, syntactic, and 084 semantic—to provide a comprehensive assessment of MLLM performance across varied complexities. 085 Moreover, we employ the rate at which questions are simultaneously answered correctly with both images as the primary metric for model evaluation. Unlike traditional visual question-answering 087 (VQA) datasets, which often rely on direct visual cues for answering questions, our benchmark 088 demands a deeper understanding of the intricate interplay between textual content and visual context. 089 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.

091 Overall, our main contributions are as follows: (a). Introduction of MMA Benchmark: We establish 092 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 094 of model evaluation. (b). Comprehensive Model Evaluation: Initial assessments of 16 MLLMs 095 reveal a significant discrepancy between model and human performances, with models averaging 096 53.22% 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-097 specific visual information. (c). Analysis of Ambiguity Types: Across the types of ambiguity, models 098 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 100 **Model Types:** A comparative analysis indicates that open-sourced MLLMs generally underperform 101 compared to proprietary MLLMs by approximately 12.59%, with Claude 3.5 Sonnet leading at 102 74.32% accuracy. 103

104 105

106

RELATED WORK 2

Multimodal large language models Recent advancements in MLLMs have opened new avenues 107 for addressing complex interaction understanding by leveraging the integration of textual and visual

2

069

070

Table 1: Comparison of different datasets with a focus on ambiguity, where Lexical/Syntactic and Semantic denote the ambiguity type.

10	Dataset	Modalities	Reasoning	Question Type	Task Type	Lexical	Syntactic	Semantic
14	WiC (Pilehvar & Camacho-Collados, 2019)	Text	×	Classification	Word Sense Disambiguation	1	×	×
	CoNLL-2012 (Pradhan et al., 2012) SemEval-2018 Task 7 (Buscaldi et al., 2017)	Text Text	×	Coreference Resolution Similarity Scoring	Coreference Resolution Semantic Similarity	×	×	×
12	AmbiEnt (Liu et al., 2023a) AmbigQA (Min et al., 2020a)	Text Text	×	Natural Language Inference QA	Ambiguity Identification Ambiguity Verification	×	×	1
3	AmbigMT (Pilault et al., 2023) AmbiCoref (Yuan et al., 2023a)	Text Text	×	MT Quality Coreference Ouality	Ambiguity in Translation Coreference Ambiguity	1	×	×
4	LAVA (Berzak et al., 2015) MM-Star (Chen et al., 2024a)	Images, Text	1	Matching Multiple Choice	Visual and Language Ambiguity Multi-task	×	×	×
5	UNPIE (Chung et al., 2024)	Images, Text	1	Open-ended	Grounding, Disambiguation, Reconstruction	1	×	×
6	MMA (Our Dataset)	Images, Text	1	Multiple Choice	Visual Question Answering	1	1	
	data. Early research, suc	ch as LXM	ERT (T	an & Bansal, 20	19), UNITER (Chen et	al., 2	2020), "	VinVL
	(Zhang et al., 2021). Vi	LBERT (Li	ı et al	2019), and VLF	P (Chen et al., 2023), fo	cuse	d on cr	eating
, ,	ioint representations to	improve n	nodality	v synergy, utiliz	ing pre-trained visual	repre	sentati	ons to
	minimize training comp	lexity. Mo	re recer	nt models, includ	ding CLIP (Radford et a	al., 20)21). A	LIGN
	(Li et al 2021) SimVl	M (Wang	et al '	2022) CoCa (Y	'u et al 2022) Flamin	σο (A	Alavrac	et al
2	(2022) BLIP-2 (Li et al	1 2023	nstruct]	BLIP2 (Lietal	2023) Mini-GPT-4 (Zhu	et al	2023)
5	Intern VI (Chen et al	2024h	WEN V	VI (Bai et al. 2	(0.23) and LL aVA (Liu	et al	2024) have
	trained visual represent	202+0), Q	VILIN-	om scratch with	massive amounts of w	oh da	., 2027 to och	ioving
	significant success in V	10000 using	ontioni	ng taske Ugwa	ver current evoluation	e moi	na, ach	
	basic visual tasks and h	VA allu Ca	lequate	ly addressed by	ndling ambiguous inn	s mai	niy 100	Pacant
	banchmarka lika 2 A M (Mo at al	0024 V	$VIS \Lambda (1 + ot of 1)$	2022) and MMMIT	n que	$-1 - 20^{\circ}$	(2) and
	beginning to incorporate	more com	2024), V Inlex an	v ISA (LI et al., A	2022, and MINIMU (It enarios into their evolution	ation	$a_{1.}, 20_{2}$	ols
	Visual question on an	ing Sime	o the in	traduction of the	a Viewal Ouastian Arrest	unonin -		(13.)
	(Antal at al. 2015) (1)	ring Sinc	e ine in	inounction of the	e visual Question Ansv	vering	g (VQA	(71)
	(Antoi et al., 2015), the	re nas been	i signifi	icant progress in	i integrating visual and			
	et al., 2010; Krishna et	al., 201/;	Goyal	et al., 2017; Hu	iuson & Manning, 201	9; L1	et al.,	2019;
	Dosovitskiy et al., 2020). However	the ch	allenge of accur	ately interpreting this c	ombi	ined da	ta still
	remains. The VQA v2 c	lataset (Go	yal et a	1., 2017) tackles	these complexities by	utiliz	ing ba	Ianced
	image pairs to enhance	detailed vi	sual an	alysis. Studies l	ike (Stengel-Eskin et a	1., 20	22) cre	ated a
	VQA dataset featuring a	mbiguous e	xample	s where images	provide just enough info	ormat	10n to a	answer
	the questions but do not	t resolve th	e inher	ent ambiguities	within the questions th	emse	elves.	Unlike
	many VQA datasets tha	t primarily	rely on	ı straightforward	l visual cues for answer	ing q	uestio	ns, our
	benchmark requires a de	eper unders	standing	g of the nuanced	interplay between text a	und vi	isual co	ontexts.
	This approach focuses o	n clarifying	g ambig	guities that arise	from the combination of	of tex	t and i	mages,
	where the contextual inf	ormation fi	rom the	images is crucia	al for disambiguating th	ie tex	tual co	ntent.
	Datasets for ambiguity	The field	d of am	biguity resolution	on in machine learning	has b	een ex	plored
	through various special	ized datase	ets, eac	h targeting spec	cific aspects of ambigu	ity.	For ex	ample.
	WiC (Pilehvar & Camac	cho-Collado	os, 2019	9) and CoNLL-2	2012 (Pradhan et al. 20	12) fo	ocus or	1 word
	sense disambiguation ar	nd coreferen	nce reso	olution, respectiv	velv, addressing text-ba	sed a	mbigu	ities in
	linguistic contexts. Data	sets like Se	emEval	-2018 Task 7 (B	uscaldi et al. 2017) A	mbiF	Ent (Li	1 et al
	2023a), AmbigOA (M	in et al 20	(20b)	AmbigMT (Pile	ault et al., 2023) and	Amhi	iCoref	(Yuan
	et al. $2023h$ further th	is work hv	tacklir	ng different form	ns of textual ambiguiti	es. fr	om sei	mantic
	similarity to natural land	mage infere	ence an	d machine transl	lation. While these data	sets o	offer v	luable
	insights they are largely	limited to	single_r	nodal text-based	tasks each focusing or	n a en	ecific t	vne of
	ambiguity The advent of	f multimod	al data	sets such as I A	VA (Berzak et al. 2015)	nasp MM	M_Star	(Chen
	et al 2024a) and MM	MII (Vue e	a_1 uatas	(0.23) represent	s significant progress h	v, v	arotin	g hoth
	visual and taxtual data	hallenging	n al., 20	(23), represent	s significant progress b	y mu	graun	g boul
	visual and textual data, (nanenging	onfinad	s to resolve and $1 to space for the start set of the space of the start set of the sta$	iguittes across modalitie	58. H(, mese
	munimodal datasets offe	thou main c	unined	i to specific tasks	s or amoiguity types. Ex	xistin	g work	s nave
	they are limited to disar	biguation v	vithin a	bus on text ambi	guides and lack multin	often	uatase	$\pi s; (2)$
	one particular type of ar	nbiguity O	viuiii S	pecific scenarios	ercome these limitation	s by i	incorpo	ss ollly prating
	multimodal data and and	noiguity. O	a wide	range of ambig	uity types to explore ar	s Uy I himi	ity icen	naung
	more general context	ompassing	a with	initige of anoig	any types to explore all	uigui	ity 1550	c5 m a
	more general context.							

158 3 BENCHMARK CONSTRUCTION

159

108

109

Our goal is to evaluate the MLLM performance under varying conditions of ambiguity. To achieve this, we introduce a comprehensive benchmark, MMA, designed to evaluate MLLM's ability to



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.

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 182 correct answer in VQA tasks.

3.1 OVERVIEW OF OUR MMA BENCHMARK

186 In order to systematically explore the capability of MLLMs to perceive and resolve ambiguities of 187 varying complexities, we categorize ambiguities into lexical, syntactic, and semantic types based on 188 the linguistic characteristics (detailed in Section 3.2). The benchmark tasks are structured as multiple-189 choice VOA scenarios, a format that simplifies the evaluation process, where **the meaning of each** 190 question is ambiguous, and they are associated with multiple images that provide varying contexts, allowing the same question to elicit different correct responses based on the visual information 191 provided as shown in Figure 2. This design forces the MLLMs to adeptly integrate and interpret both 192 textual and visual data to select the most accurate answer, reflecting the true potential and challenges 193 of deploying such models in diverse, ambiguity-filled environments. 194

3.2 **Types of Ambiguities** 196

175

176

177

178 179

181

183

185

195

197 We divide the ambiguity into the following types and design questions according to each different ambiguity type. Each category is designed to evaluate specific aspects of how well models integrate 199 and interpret complex linguistic and visual information to resolve ambiguities. For a more detailed 200 understanding, we provide examples of each type in Table 2.

201 **Lexical ambiguities.** Lexical ambiguity mainly evaluates the ambiguity caused by polysemy in 202 sentences. We considered the ambiguity caused by nouns, adjectives, and verbs. The verb category 203 includes both the ambiguity of polysemy and the ambiguity of different emotions it may evoke.

204 Syntactic ambiguities. Syntactic ambiguities occur when sentence structures allow for multiple 205 interpretations. There are three main types: (a) Attachment Ambiguity: This occurs when a 206 modifying phrase, usually a prepositional phrase or clause, can logically attach to more than one 207 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 208 the modifiers apply to all or just some elements. (c) Structural Ambiguity: This arises when verbs 209 can be used in both transitive and intransitive forms, leading to different meanings. 210

Semantic ambiguities. Semantic ambiguities involve the broader meanings of text and their 211 interaction with visual elements : (a) Idiomatic Ambiguity: This occurs with idiomatic expressions 212 that can be interpreted both literally and metaphorically. (b) **Pragmatic Ambiguity**: This arises from 213 interpreting a sentence in different contexts provided by visual cues, affecting how the listener or 214 viewer understands the relevance and expected response. 215

Example	Scenario 1	Explanation 1	Scenario 2	Explanation 2	Туре
The meaning of <u>"bat"</u> .		One type of animal	X	The tool used in baseball	Lexical (Noun)
She saw the cat <u>under the tree</u> .		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.	Coordinatio (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.	E.E.	She literally has a green-colored thumb.		She is skilled at gardening.	Idiomatic (Semantic)
Everyone is not <u>here</u> .	×××	No one is here.	***	Not everyone is here.	Pragmatic (Semantic)

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

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.

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 (Rombach et al., 2022) and Dall-E (OpenAI, 2024a). 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.

274 3.4 HUMAN EVALUATION 275

To explore how humans perform on our MMA benchmark, we invite five annotators with near-native proficiency whose English level meets the CEFR¹ 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:

- 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
- Explore the reasons why MLLMs lag behind humans in MMA benchmark? Sec 4.3.2
- How well do the models handle each type of ambiguity? Sec 4.3.4
- To what extent does model scale (number of parameters) influence performance? Sec 4.3.5





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

4.1 EVALUATION MODELS

Figure 4: The ablation study about the parameter number and the ambiguity accuracy performance on different ambiguity types.

We evaluate 17 recent multimodal LLMs on our benchmark, including 6 proprietary MLLMs: GPT-4-vision (OpenAI, 2024b), GPT-40 (OpenAI, 2024c), Claude-3-Opus (Anthropic, 2024), Gemini-1.5-Pro (DeeoMind, 2024), Gemini-1.0-Pro-Vision (DeeoMind, 2023), Claude 3.5 Sonnet (Anthropic, 2024) and 11 open-source MLLMs: LLaVAV-Next (Liu et al., 2024), VILA1.5 (Lin et al., 2023), Yi-VL-34b (AI et al., 2024), InternVL-Chat-V1.5(Chen et al., 2024b), InternVL2(Chen et al., 2024b), CogVLM2-LLaMa3-Chat-19B (Wang et al., 2023), DeepSeek-VL-7b-Chat (Lu et al., 2024), MiniCPM-Llama3-V-2.5(OpenBMB, 2024), HPT1.5-Air (HYPERGAI, 2024), Qwen2-V(Wang et al., 2024b), LLaVA-OneVision(Li et al., 2024). Our evaluation is conducted under a zero-shot setting.Due to the page limit, we describe these models in detail in Appendix.

4.2 EVALUATION METRICS

Evaluating the ability of MLLMs to resolve ambiguity in multimodal settings requires metrics that
 go beyond standard accuracy measures. To capture the nuances of this challenge, we introduce this
 metrics for the MMA benchmark:

¹https://www.cambridgeenglish.org/exams-and-tests/cefr/

Ambiguity Accuracy (Amb_A) This metric is calculated as the percentage of questions where the
 model correctly answer for both paired images. A high Amb_A indicates that the model does not
 simply latch onto one possible interpretation of the ambiguity. Instead, it effectively integrates visual
 information from images to arrive at the most appropriate answer for each scenario. The examples
 are given in Figure 2.

4.3 MAIN RESULTS

329 330

333

351

352

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

		Adjective (30)	Noun (238)	Verb (16)	Attachment (24)	Coordination (46)	Structural (14)	Pragmatic (132)	Idiom (22)	Lexical (284)	Syntactic (84)	Semantic (154)	Overall (522)
					Prop	rietary MLLMs:	:						
	GPT-4 Vision (OpenAI, 2024b)	0.87	0.748	0.63	0.23	0.41	0.29	0.68	0.62	0.75	0.33	0.65	0.65
	GPT-40-0513 (OpenAI, 2024c)	0.80	0.82	0.88	0.08	0.41	0.43	0.65	0.73	0.82	0.31	0.69	0.70
	Gemini 1.5 Pro(DeeoMind, 2024)	0.79	0.75	0.83	0.54	0.59	0.14	0.74	0.38	0.76	0.50	0.57	0.66
	Gemini 1.0 Pro Vision (DeeoMind, 2023)	0.69	0.68	0.40	0.00	0.32	0.00	0.41	0.29	0.67	0.17	0.35	0.49
	Claude 3 Opus (Anthropic, 2024)	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 (Anthropic, 2024)	0.83	0.83	0.86	0.67	0.47	0.50	0.73	0.30	0.83	0.53	0.67	0.74
	GPT-40-0806 (OpenAI, 2024c)	0.80	0.82	0.75	0.31	0.55	0.57	0.73	0.65	0.81	0.48	0.69	0.72
	Proprietary Average	0.79	0.75	0.67	0.26	0.41	0.28	0.60	0.45	0.75	0.34	0.55	0.62
					Open	-source MLLMs	:						
	LLaVA-NeXT-34B (Liu et al., 2024)	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 (Liu et al., 2024)	0.67	0.64	0.38	0	0.09	0	0.33	0.59	0.63	0.05	0.45	0.48
LLaV	LLaVA-NeXT-7B (Liu et al., 2024)	0.13	0.60	0.13	0	0	0.14	0.28	0.27	0.52	0.02	0.27	0.37
	VILA1.5-40b (Lin et al., 2023)	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 (Lin et al., 2023)	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 (Lin et al., 2023)	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 (AI et al., 2024)	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 (Chen et al., 2024b)	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 (Chen et al., 2024b)	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 (Wang et al., 2023)	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 (Lu et al., 2024)	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 (OpenBMB, 2024)	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 (HYPERGAI, 2024)	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 (Wang et al., 2024b)	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 (Wang et al., 2024b)	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 (Li et al., 2024)	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 Li et al. (2024)	0.47	0.74	0.38	0.23	0.45	0.00	0.50	0.18	0.69	0.31	0.45	0.56
ļ	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

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 Claude 3.5 Sonnet (74%), achieve the best overall performance on
 Amb_A and excel at handling lexical ambiguities (83%). Among open-source models, InternVL Chat-V1-5 (69.7%) shows strong performance, particularly in lexical categories (82%), achieving
 nearly comparable performance to Claude 3.5 Sonnet.

361 Despite these advancements, the best-performing models like Claude 3.5 Sonnet and GPT-40 still 362 show a substantial gap when compared to human performance. Claude 3.5 Sonnet achieves an 363 overall accuracy of 74%, which is 15% lower than the human benchmark of 89%. Similarly, GPT-40 364 performs 19% lower than human performance with an overall accuracy of 70%. Gemini-1.5 pro and InternVL-Chat-V1-5 also underperform humans by 23% and 19%, respectively, with overall accuracy of 66% and 70%. This significant performance gap is particularly evident in tasks involving 366 syntactic and semantic ambiguities. For example, Claude 3.5 Sonnet and GPT-40 achieve accuracy 367 of 53% and 31% in syntactic ambiguities, respectively, compared to the human accuracy of 89%. 368 In semantic ambiguities, Claude 3.5 Sonnet and GPT-40 achieve 67% and 69%, respectively, while 369 humans achieve 85%. 370

- 371
- 372

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 assessed the accuracy of MLLMs when they were provided solely with text inputs. The metric

Table 4: MLLM Performance with Text-Only Input: We assessed the ratio of selecting one of
 the correct answers when MLLMs are given text-only input. This metric is used to measure the
 language understanding ability of MLLMs, addressing concerns about the potential bias introduced
 by 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-40-2024-05-13	0.85	0.89
GPT-40-2024-08-06	0.85	0.88
InternVL-Chat-V1-5	0.85	0.86

Table 5: MLLMs' Error Consistency Rate: This metric represents the ratio of instances where
 MLLMs provide the same answer even when presented with two different images. It is used to
 measure the extent to which MLLMs neglect image information in clarifying ambiguities during the
 question-answering process.

Model	Lexical	Syntactic	Semantic	Overall
Claude 3 Opus	0.86	0.72	0.89	0.84
GPT-4o-2024-05-13	0.72	0.83	0.79	0.78
InternVL-Chat-V1-5	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

400 401

394

397

399

382

402 used represents the rate at which the model's response matches one of the correct answers in each 403 pair of data (ambiguity pair), it is considered accurate. As shown in Table 4, MLLMs demonstrate 404 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 405 strong across lexical, syntactic, and semantic categories, with most models scoring above 80% in 406 each. Claude 3.5 Sonnet shows the most balanced performance across categories, while others like 407 InternVL-Chat-V1-5 exhibit some variability (e.g., 90% lexical vs. 74% syntactic). These results 408 indicate that minor textual issues have minimal impact on MLLMs' ability to select correct answers. 409

MLLMs' Error Consistency Rate: This Error Consistency Rate (ECR) - defined as the rate of 410 selecting the same answer among incorrect cases. As shown in Table 5, when MLLMs made errors, 411 they demonstrated a high consistency rate in choosing the same option twice. This rate ranged from 412 71% to 84% overall, depending on the model. The consistently high rates across lexical, syntactic, 413 and semantic levels indicate that these models often failed to effectively leverage visual information 414 when answering questions. Instead, they exhibited a strong bias towards the text modality, relying 415 primarily on textual cues even when visual information was available. More error analysis are given 416 in Appendix A.5. 417

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.

424 425

426

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 57.70% in Amb_A, while open-sourced models
 obtain 47.04% in Amb_A as Table 3 shows. For both indicators, proprietary models outperform
 open-sourced models.

1.0 Ambiguity Accuracy Syntactic 0.8 Semantic 0.6 0.4 0.2 Laude 3.5 somet LIANA ORENEONTRE Spt. Avision preview 0emini-1.5 claude 3 vision 18484-ret-349 Over2 VL 128 optomay VILAADO Internal 1.5 0.0 Cogviniz DeepSeetVI HPT AIR VILA 130 optaoAug

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

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-40 shows 82% accuracy on lexical ambiguities but drops to 31% and 69% on syntactic and semantic ambiguities, respectively.

456 Syntactic ambiguities present a unique challenge because they involve the relationships between 457 components within a sentence. Often, even a short modifier can introduce ambiguity, making it 458 difficult for models to resolve these cases without fine-grained analysis. To effectively handle 459 syntactic ambiguities, models need not only a more granular approach to language processing but also 460 the capability to accurately recognize positional relationships and details in images. This requires a 461 higher level of precision compared to lexical and semantic ambiguities. Models like Cogvlm2 and 462 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%).

466 467

432

433

434

435

436 437

438

444

445 446

447

448

4.3.5 SCALING LAW ON MMA

468 To comprehend whether the parameter number affects performance on the MMA benchmark, we 469 conducted experiments on the same series of models with varying sizes, all trained on similar data. 470 As Figure 4 shows, there is a clear improvement in ambiguity accuracy as the parameter count 471 increases across different ambiguity types. For instance, larger models like VILA1.5-40B consistently 472 outperform smaller ones such as VILA1.5-3B. The larger models show significant improvements 473 in handling lexical, syntactic, and semantic ambiguities, demonstrating that increased parameters enhance the model's ability to understand and disambiguate complex, multi-modal inputs. This trend 474 indicates a positive correlation between model size and performance on the MMA benchmark. 475

476 477

478

5 LIMITATION

479 Data collection Due to constraints on the number of participants, the dataset size is limited in 480 certain categories. Despite this limitation, we emphasize that the quality and representativeness of 481 the dataset are more crucial for establishing a meaningful benchmark than merely the number of 482 samples. As demonstrated in Table 3, the considerable performance discrepancy between human 483 participants and MLLM responses underscores the benchmark's effectiveness in highlighting the 484 current challenges that MLLMs face, particularly their inability to adequately utilize visual context to 485 resolve textual ambiguities. Moving forward, we are committed to expanding the dataset in future 486 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. 494 The images in our benchmark are specifically chosen to provide the necessary context to clarify 495 ambiguities in the accompanying texts. Due to the current limitations of search engines, which 496 struggle with semantic search, it is challenging to find suitable images that naturally align with the 497 required context (This doesn't mean that these images don't exist.). Therefore, using generated 498 images is the most effective approach. They are instrumental in simulating the diverse and often 499 unconventional situations that MLLMs encounter in real applications. MLLMs are expected to 500 perform comparably to humans in these scenarios, regardless of the variability in inputs. However, 501 our human study shows that humans can achieve approximately 90% accuracy on this benchmark 502 without any additional interactions. This sharply contrasts with the average accuracies of 58% for 503 closed-source models and 47% for open-source models.

504 505

506

493

6 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.

512

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.

526 527 528

7 CONCLUSION

529 This paper introduces MMA, the first benchmark designed specifically to evaluate the ability of 530 Multimodal Large Language Models (MLLMs) to understand and respond to ambiguous queries. 531 MMA leverages a multiple-choice visual question-answering format, presenting MLLMs with a 532 question and two images depicting contrasting scenarios that lead to different correct answers. 533 Our evaluation of 16 MLLMs, including both limited-access and open-sourced models, reveals a 534 significant performance gap compared to human performance. While humans achieve an accuracy 535 of 88.97%, the MLLMs average only 50.59% accuracy. This indicates a fundamental challenge for current MLLMs: effectively integrating scenario-specific visual information to disambiguate 536 questions and arrive at the correct answer. Even the top-performing model, GPT-40 and Claude3.5-537 Sonnet, attains only about 70.00% accuracy, highlighting considerable room for developing MLLMs 538 that can effectively leverage visual information to clarify the textual ambiguity and capable of human-level understanding and reasoning in complex, real-world scenarios.

540 REFERENCES

555

- 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.
- 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.
- 553 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.
- 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*, pp. 2425–2433, 2015.
- Soufia Bahmani. Enhancing human-computer interaction through large language models: Opportunities, challenges, and future directions.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023.
- 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 (eds.), *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 1477–1487, Lisbon, Portugal, September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1172. URL https://aclanthology.org/D15-1172.
- 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)*, pp. 679–688, 2017.
- Fei-Long Chen, Du-Zhen Zhang, Ming-Lun Han, Xiu-Yi Chen, Jing Shi, Shuang Xu, and Bo Xu.
 Vlp: A survey on vision-language pre-training. *Machine Intelligence Research*, 20(1):38–56, 2023.
- 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?, 2024a.
- 579
 580
 581
 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.
- 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, 2024b.
- 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.
- 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.

594 595	Google DeeoMind. Gemini pro vision, 2023. URL https://deepmind.google/ technologies/gemini/pro/.										
597 598	Google DeeoMind. Gemini pro1.5, 2024. URL https://deepmind.google/ technologies/gemini/pro/.										
599 600 601 602	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. <i>arXiv preprint arXiv:2010.11929</i> , 2020.										
603 604 605 606	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 <i>Extended Abstracts of the CHI Conference on Human Factors in Computing Systems</i> , pp. 1–11, 2024.										
608 609 610	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 <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 6904–6913, 2017.										
611 612 613	Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 6700–6709, 2019.										
614 615 616 617	HYPERGAI.Hpt1.5air:Bestopen-sourced8bmultimodalllmwithllama3,2024.URLhttps://www.hypergai.com/blog/hpt-1-5-air-best-open-sourced-8b-multimodal-llm-with-llama-3.										
618 619	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.										
620 621 622 623	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. <i>International journal of computer vision</i> , 123:32–73, 2017.										
624 625 626	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. <i>arXiv preprint arXiv:2408.03326</i> , 2024.										
627 628 629 630	Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. <i>Advances in neural information processing systems</i> , 34:9694–9705, 2021.										
631 632	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.										
633 634 635	Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. <i>arXiv preprint arXiv:1908.03557</i> , 2019.										
636 637 638	Yihang Li, Shuichiro Shimizu, Weiqi Gu, Chenhui Chu, and Sadao Kurohashi. Visa: An ambiguous subtitles dataset for visual scene-aware machine translation. <i>arXiv preprint arXiv:2201.08054</i> , 2022.										
639 640	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.										
642 643	Jieru Lin, Danqing Huang, Tiejun Zhao, Dechen Zhan, and Chin-Yew Lin. Designprobe: A graphic design benchmark for multimodal large language models. <i>arXiv preprint arXiv:2404.14801</i> , 2024.										
644 645 646 647	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 (eds.), <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , December 2023a. URL https://aclanthology.org/2023.emnlp-main.51.										

648 Alisa Liu, Zhaofeng Wu, Julian Michael, Alane Suhr, Peter West, Alexander Koller, Swabha 649 Swayamdipta, Noah A Smith, and Yejin Choi. We're afraid language models aren't modeling 650 ambiguity. arXiv preprint arXiv:2304.14399, 2023b. 651 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 652 Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL https:// 653 llava-vl.github.io/blog/2024-01-30-llava-next/. 654 655 Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, 656 Zhuoshu Li, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. Deepseek-vl: 657 Towards real-world vision-language understanding, 2024. 658 Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic 659 representations for vision-and-language tasks. Advances in neural information processing systems, 660 32, 2019. 661 662 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 663 Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language 664 Resources and Evaluation (LREC-COLING 2024), Torino, Italia, May 2024. ELRA and ICCL. 665 URL https://aclanthology.org/2024.lrec-main.1. 666 667 Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. AmbigQA: Answering 668 ambiguous open-domain questions. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu 669 (eds.), Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing 670 (EMNLP), November 2020a. URL https://aclanthology.org/2020.emnlp-main. 671 466. 672 Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. Ambigga: Answering 673 ambiguous open-domain questions. arXiv preprint arXiv:2004.10645, 2020b. 674 675 Sigrid Norris. Analyzing multimodal interaction: A methodological framework. Routledge, 2004. 676 OpenAI. Dall-e-3, 2024a. URL https://openai.com/index/dall-e-3/. 677 678 OpenAI. Gpt-4 technical report, 2024b. 679 **OpenAI.** Hello gpt-40, 2024c. URL https://openai.com/index/hello-gpt-40/. 680 681 OpenBMB. Minicpm-llama3-v 2.5, 2024. URL https://github.com/OpenBMB/ 682 MiniCPM-V?tab=readme-ov-file#minicpm-llama3-v-25. 683 Jonathan Pilault, Xavier Garcia, Arthur Bražinskas, and Orhan Firat. Interactive-chain-prompting: 684 Ambiguity resolution for crosslingual conditional generation with interaction. In Jong C. 685 Park, Yuki Arase, Baotian Hu, Wei Lu, Derry Wijaya, Ayu Purwarianti, and Adila Alfa Kris-686 nadhi (eds.), Proceedings of the 13th International Joint Conference on Natural Language 687 Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Com-688 putational Linguistics (Volume 1: Long Papers), pp. 455-483, Nusa Dua, Bali, November 689 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.ijcnlp-main.31. URL 690 https://aclanthology.org/2023.ijcnlp-main.31. 691 Mohammad Taher Pilehvar and Jose Camacho-Collados. WiC: the word-in-context dataset for 692 evaluating context-sensitive meaning representations. In Proceedings of the 2019 Conference of 693 the North American Chapter of the Association for Computational Linguistics: Human Language 694 Technologies, Volume 1 (Long and Short Papers), Minneapolis, Minnesota, June 2019. Association 695 for Computational Linguistics. URL https://aclanthology.org/N19-1128. 696 697 Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. Conll-2012 shared task: Modeling multilingual unrestricted coreference in ontonotes. In Joint conference 698 on EMNLP and CoNLL-shared task, pp. 1-40, 2012. 699 700

702 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, 703 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual 704 models from natural language supervision. In International conference on machine learning, pp. 705 8748-8763. PMLR, 2021. 706 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. Highresolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF confer-708 ence on computer vision and pattern recognition, pp. 10684–10695, 2022. 709 710 Elias Stengel-Eskin, Jimena Guallar-Blasco, Yi Zhou, and Benjamin Van Durme. Why did the 711 chicken cross the road? rephrasing and analyzing ambiguous questions in vqa. arXiv preprint arXiv:2211.07516, 2022. 712 713 Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from trans-714 formers. arXiv preprint arXiv:1908.07490, 2019. 715 716 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 717 multimodal models. arXiv preprint arXiv:2312.11805, 2023. 718 719 Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J Zico Kolter, Louis-Philippe Morency, and 720 Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In 721 Proceedings of the conference. Association for computational linguistics. Meeting, volume 2019, 722 pp. 6558. NIH Public Access, 2019. 723 Matthew Turk. Multimodal interaction: A review. Pattern recognition letters, 36:189–195, 2014. 724 725 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz 726 Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing 727 systems, 30, 2017. 728 Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao 729 Sang. Mobile-agent: Autonomous multi-modal mobile device agent with visual perception. arXiv 730 preprint arXiv:2401.16158, 2024a. 731 732 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, 733 Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the 734 world at any resolution. arXiv preprint arXiv:2409.12191, 2024b. 735 Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, 736 Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. 737 Cogvlm: Visual expert for pretrained language models, 2023. 738 739 Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. Simvlm: Simple visual language model pretraining with weak supervision, 2022. 740 741 Thomas Wasow, Amy Perfors, and David Beaver. The puzzle of ambiguity. Morphology and the web 742 of grammar: Essays in memory of Steven G. Lapointe, pp. 265–282, 2005. 743 744 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: 745 10.1109/TPAMI.2023.3275156. 746 747 Yi Yang, Qingwen Zhang, Ci Li, Daniel Simões Marta, Nazre Batool, and John Folkesson. Human-748 centric autonomous systems with llms for user command reasoning. In Proceedings of the 749 IEEE/CVF Winter Conference on Applications of Computer Vision, pp. 988–994, 2024. 750 Keen You, Haotian Zhang, Eldon Schoop, Floris Weers, Amanda Swearngin, Jeffrey Nichols, Yinfei 751 Yang, and Zhe Gan. Ferret-ui: Grounded mobile ui understanding with multimodal llms. arXiv 752 preprint arXiv:2404.05719, 2024. 753 754 Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. 755 Coca: Contrastive captioners are image-text foundation models, 2022.

756 757	Yuewei Yuan, Chaitanya Malaviya, and Mark Yatskar. AmbiCoref: Evaluating human and model sensitivity to ambiguous coreference. In Andreas Vlachos and Isabelle Augenstein (eds.), <i>Findings</i>
758	of the Association for Computational Linguistics: EACL 2023, Dubrovnik, Croatia, May 2023a.
759	Association for Computational Linguistics. URL https://aclanthology.org/2023.
760	findings-eacl.75.
761	
762 763	Yuewei Yuan, Chaitanya Malaviya, and Mark Yatskar. Ambicoref: Evaluating human and model sensitivity to ambiguous coreference. <i>arXiv preprint arXiv:2302.00762</i> , 2023b.
764	Viena Vuo Vuonshana Ni Vai Zhana Tianuu Zhana Duogi Liu Ca Zhana Samuel Stavana Donafu
765	Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin,
766	Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen.
767	Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert
768	agi, 2023.
769	Chaoyun Zhang Ligun Li, Shilin He, Yu Zhang Be Oige, Si Oin, Minghua Ma, Yu Kang, Oingwai
770	Lin Sarayan Raimohan et al. Ufo: A ui-focused agent for windows os interaction arXiv preprint
771	arXiv:2402.07939. 2024.
772	
773	Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and
774	Jianfeng Gao. Vinvl: Revisiting visual representations in vision-language models, 2021.
775	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing
776	vision-language understanding with advanced large language models, 2023.
777	Vuka Zhu, Olivar Groth, Michael Barnstein, and Li Fai Fai, Visual 7111; Grounded question answering
778 779	in images. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp.
780	4995–5004, 2016.
781	
782	
783	
784	
785	
786	
787	
788	
789	
790	
791	
792	
793	
794	
705	
795	
795 796	
796 797	

810 A APPENDIX

831 832

833 834

835

836 837 838

839

840

841

842

843

844 845

846

812 A.1 DISTRIBUTION OF DATASET

As shown in Figure 6, the MMA dataset consists of 522 images and 261 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).



Figure 6: Ambiguity Type Composition of MMA benchmark

A.2 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.

A.3 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.

A.4 ABLATION STUDY

847 Same images with lexical or semantic questions To understand why MLLMs perform better on 848 lexical ambiguity compared to semantic ambiguity, we explored how changing the question type on 849 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 850 reasoning into the question. For example, given an image of a table, a synonym question for lexical 851 ambiguity might be "What is the meaning of table?" where the model identifies "table" as a piece 852 of furniture. In contrast, a reasoning question for semantic ambiguity would be "How can we best 853 utilize the space on this table?" which requires the model to consider various uses of the table. This 854 type of question tests the model's ability to perform object grounding and higher-order reasoning, 855 areas where MLLMs often show weaker performance due to their reliance on pattern recognition 856 rather than true comprehension. More examples are given in Appendix. 857

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.

A.5 ERROR ANALYSIS

879

880

882

883

884

885

887 888

889

890

891

892

893 894 895

896

897

902

903

909

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-40 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.



Figure 8: Error type distribution of GPT-40,
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).

P15
 P16
 P17
 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



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.

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.

A.6 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-

	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 6: Five people have different performance across different types of ambiguities



1012 1013

1023

1025



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.

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 answer options with subtle differences, or struggling to correlate images with text (particularly when the text contains advanced expressions or extended meanings). Here is an example where all respondents have failed to provide the correct answer Figure 13.





1014 A.7 SCALING LAW WITH LLAVA SERIES MODELS

As for Llava series models, the scaling law holds true for most metrics, with performance continuously improving as the model's parameter count increases(see Figure 9). This provides further evidence for the scaling law on the MMA benchmark. However, there is an exception when it comes to semantic ambiguity, where the middle-sized model performs best.

1020 A.8 CASE STUDY

1022 More examples of GPT-40 will be presented in this section.

1024 Example of GPT-40

1. Example of Coordination Ambiguity





Figure 18: Example of Attachment Ambiguity



Figure 21: Example of Verb Ambiguity