MUCAR: Benchmarking Multilingual Cross-Modal Ambiguity Resolution for Multimodal Large Language Models

Anonymous ACL submission

Abstract

Multimodal Large Language Models (MLLMs) have demonstrated significant advances across numerous vision-language tasks. Due to 004 their strong image-text alignment capability, MLLMs can effectively understand image-text pairs with clear meanings. However, effec-007 tively resolving the inherent ambiguities in natural language and visual contexts remains challenging. Existing multimodal benchmarks typically overlook linguistic and visual ambiguities, relying mainly on unimodal context for disambiguation and thus failing to exploit the mutual clarification potential between modalities. To bridge this gap, we intro-015 duce MUCAR, a novel and challenging benchmark designed explicitly for evaluating mul-017 timodal ambiguity resolution across multilingual and cross-modal scenarios. MUCAR includes: (1) a multilingual dataset where ambiguous textual expressions are uniquely resolved by corresponding visual contexts, and (2) a dual-ambiguity dataset that systematically pairs ambiguous images with ambiguous textual contexts, with each combination carefully constructed to yield a single, clear interpretation through mutual disambiguation. Extensive evaluations involving 19 state-of-the-art multimodal models-encompassing both opensource and proprietary architectures-reveal substantial gaps compared to human-level performance, highlighting the need for future research into more sophisticated cross-modal ambiguity comprehension methods, further pushing the boundaries of multimodal reasoning.

1 Introduction

037

041

Multimodal Large Language Models (MLLMs; OpenAI 2023, 2024; Liu et al. 2023a; Dai et al. 2023) have advanced significantly in handling diverse inputs like text and images, with contextual understanding being key to their success in tasks such as question answering (Shao et al., 2023; Liu et al., 2023b; Antol et al., 2015), image captioning (Luo et al., 2023; Wang et al., 2023; Chen



Figure 1: The interpretation of ambiguous text and visuals can be shaped by context and scenario. For instance, in "My wife and I met a couple on the boat," it is unclear whether "on the boat" modifies "met" or "a couple," while the image leaves the identity of the wife uncertain. Varying the visual setting (e.g., riverbank vs. cruise deck) and textual cues (e.g., age of the speaker) yields different resolutions. Each of the four context-image combinations leads to a distinct answer, with reasoning color-coded: purple (Context 1), red (Context 2).

et al., 2015), and multimodal reasoning (Gupta and Kembhavi, 2023; Chen et al., 2023; Zellers et al., 2019). As MLLMs evolve, the ability to integrate multimodal context becomes crucial for accurate responses, underscoring the need for thorough evaluation of their contextual comprehension in real-world settings.

Prior studies have largely emphasized tasks with clear and unambiguous inputs (Fu et al., 2023a; Ying et al., 2024; Li et al., 2023), frequently neglecting the ambiguity that naturally arises in both visual and textual modalities. Consider the example in Figure 2, neither the context nor the image alone can resolve the question "Who is my wife, the woman on the left or the right?", where a single sentence or image can often support multiple plausible interpretations depending on the specific scenario or provided context. Figure 2 illustrates this challenge clearly. Consider the sentence "My

Benchmark	Visual Ambiguity	Context Ambiguity	Multi-Languages	Evaluator
MME (Fu et al., 2023a)	×	×	×	Metrics
MMBench (Liu et al., 2023c)	×	×	×	GPT
MMT-Bench (Ying et al., 2024)	×	×	×	GPT
MMStar (Chen et al., 2024a)	×	×	×	Metrics
HallusionBench (Guan et al., 2023)	✓	×	×	Metrics
CODIS (Luo et al., 2024)	\checkmark	×	×	Human / GPT
Illusory VQA (Rostamkhani et al., 2024)	\checkmark	×	×	Human
MHaluBench (Chen et al., 2024b)	\checkmark	×	×	GPT
MMA (Wang et al., 2024a)	×	\checkmark	×	Metrics
VAGUE (Nam et al., 2024)	×	\checkmark	×	Metrics
3AM (Ma et al., 2024)	×	\checkmark	\checkmark	Metrics
UNPIE (Chung et al., 2024)	×	\checkmark	\checkmark	GPT
MUCAR (Ours)	\checkmark	\checkmark	\checkmark	Human / GPT

Table 1: Comparison of our proposed MUCAR with recent vision-language benchmarks.

wife and I met a couple on the boat." This sentence contains structural ambiguity: it is unclear whether "on the boat" modifies the verb "met" (indicating the location of the meeting) or the noun phrase "a couple" (specifying the location of the couple). Simultaneously, the accompanying image introduces visual ambiguity concerning the referent of "the wife" among the depicted women. Notably, the ambiguity cannot be resolved independently within either modality; instead, mutual disambiguation arises when different textual scenarios (e.g., differing speaker ages) combine with different visual contexts (riverbank vs. cruise deck scenarios). Each unique combination yields a single, unambiguous interpretation, showing that textual and visual ambiguities can mutually clarify each other.

063

065

067

068

073

077

To systematically evaluate the capabilities of MLLMs to resolve such complex multimodal ambiguities, we introduce MUCAR, a novel bench-081 mark specifically designed for MUltilingual Crossmodal Ambiguity Resolution. Table 1 summarizes 083 recent benchmarks designed to evaluate MLLMs in terms of visual ambiguity, contextual ambiguity, multilinguality, and the type of evaluator used (e.g., metrics, GPT, or human annotations). While 087 early benchmarks such as MMT-Bench (Ying et al., 2024), MMStar (Liu et al., 2023c), and MME-RealWorld (Li et al., 2023) focus on general multimodal tasks, they lack coverage of ambiguityrelated phenomena. More recent benchmarks like HallusionBench (Guan et al., 2023), Illusory VQA (Rostamkhani et al., 2024), and CODIS (Luo et al., 2024) begin to explore visual ambiguity, but

often overlook contextual disambiguation or multilingual diversity. Notably, only a few benchmarks incorporate human evaluation, which is essential for assessing ambiguity understanding. To the best of our knowledge, MUCAR is the first benchmark to comprehensively address visual ambiguity, contextual ambiguity, and multilinguality, while integrating both human and GPT-based evaluation. This design enables a more rigorous and realistic assessment of ambiguity resolution capabilities in multimodal large language models.

097

098

100

101

102

103

104

105

106

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

To summarize, our main contributions are:

- We construct **MUCAR**, the first *multilingual cross-modal ambiguity resolution benchmark*, featuring 1278 manually curated samples in Chinese, English, and Malay, including uniquely designed dual-ambiguity cases.
- We systemically evaluate **19 sota MLLMs** (both open-source and closed-source), revealing significant limitations in resolving multi-lingual multimodal ambiguities.
- We propose a **simple yet effective agentbased framework** for multimodal disambiguation, which improves performance through explicit cross-modal reasoning.

2 Related Work

Context Ambiguity Ambiguity is an inherent characteristic of linguistic text, emerging naturally due to the potential for multiple interpretations, especially in open-domain question answering tasks (Min et al., 2020; Sun et al., 2023). Existing research on ambiguity resolution in language



Figure 2: Taxonomy of our benchmark. We present one example for each category. Each example includes a context (C), a question (Q), and two different images with their corresponding answers and explanations.

primarily follows two directions. One line leverages contextual cues from surrounding text to resolve ambiguity, as seen in Gao et al. (2024); Lee et al. (2025), which uses in-context learning to disambiguate textual inputs. Another line exploits visual information to disambiguate language, particularly for polysemous word translation, such as in 3AM (Ma et al., 2024), which aligns ambiguous words with visual semantics.

129

130

132

133

135

136

137

138

141

142

143

144

145

Aligned with the multimodal ambiguity setting in MMA (Wang et al., 2024a), our work focuses on using visual input to disambiguate multilingual expressions. Beyond standard NLP ambiguities, we also address those arising from domain-specific and cultural differences. In particular, we construct a benchmark featuring both textual and visual ambiguity across multiple languages, aiming to evaluate the ability of MLLMs to resolve complex multimodal ambiguities in realistic scenarios.

Visual Ambiguity Visual ambiguity often stems from incomplete visual cues or interfering noise in the scene (Denison et al., 2018). Most previous vision-language benchmarks assume unambiguous (Liu et al., 2023c; Fu et al., 2023a; Liu et al., 2023c; Li et al., 2023) input or highlight the visual ambiguities caused by optical illusions (Guan et al., 2023; Rostamkhani et al., 2024; Cui et al., 2023; Fu et al., 2023b). Early Multimodal datasets like MS-COCO (Chen et al., 2015) focus on literal descriptions, while later works, e.g., CODIS (Luo et al., 2024) highlight the need for diverse context to reflect multiple valid interpretations. Inspired by CODIS (Luo et al., 2024), we assess the capability of MLLMs to disambiguate visual ambiguity through textual modalities instead of just recognizing ambiguities. Different from CODIS,

146

147

148

149

150

151

152

153

154

155

156

157

159

160

161

162

we construct challenging dual-ambiguity instances,
combining ambiguous visuals and texts that jointly
resolve into a single interpretation, further testing
the limits of multimodal reasoning.

168 **3 MUCAR**

169

170

171

172

173

174

175

176

MUCAR is proposed for evaluating the capabilities of MLLMs in image-dependent context disambiguation. Figure 2 presents several examples from our benchmark, highlighting the diversity of contexts covered. In this section, we first describe our taxonomy of context. Then, we delve into the instruction design. Finally, we introduce data collection procedures.

3.1 Taxonomy

Given the extensive and varied nature of context 178 information, comprehensive cataloging of all forms of context is challenging. With the aim of establish-180 ing an outstanding benchmark for disambiguation, 181 we identified seven representative types. The first 182 six types are inspired by the information people 183 require to understand context. When collecting the 184 data, we surprisingly found that in some cases, the 185 combination of ambiguous text and ambiguous images led to mutual disambiguation. In other words, neither modality alone provided sufficient clarity, 188 but together they resolved the ambiguity inherent in 189 both. This observation inspired us to define the sev-190 enth type and collect the corresponding data. To the 191 best of our knowledge, we are the first to construct 192 data for this type. Figure 2 illustrates examples 193 with corresponding classification explanations. 194

Polysemy. Navigli (2009) defines the term, polysemy, by referring to a word with two or more
related meanings. These related meanings often
share a conceptual core, with one meaning typically being an extension or variation of the other.
Some contexts can be interpreted both literally and
metaphorically. The "Polysemy" part in Figure 2
gives a good example.

Homonymy. Homonymy is also defined in Navigli
(2009). Opposite to polysemy, homonymy refers
to a word having two or more unrelated meanings
that stem from different historical origins, and the
meanings of homonyms have no inherent connection. Disambiguating homonyms relies on the other
elements. The "Homonymy" part in Figure 2 serves
as a good example.

Grammar. This ambiguity occurs when sentencestructures allow for multiple interpretations, often

due to the placement of words or phrases. Such structural issues can make it unclear which part of the sentence a modifier applies to or the relationship between different clauses. The "Grammar" part in Figure 2 provides a clear illustration. As the "Grammar" part in the figure shows, in Chinese, "我的门没有锁" can be interpreted as "My door does not have a lock" or "My door has not been locked", here "锁" can be understood as a noun or verb, which leads to ambiguity. 213

214

215

216

217

218

219

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

Semantics. Understanding the timing and sequence of events is crucial when we understand a context. However, an isolated context can only provide us with static information, which is insufficient for dynamic events. The disambiguation can only be achieved when unambiguous image gives us more information. The "Semantics" case in Figure 2 provides a representative example.

Specialized. Context can be interpreted differently depending on the circumstances of the viewer, including their background knowledge, expertise, or personal situation. What is clear and unambiguous to someone with specialized training might be confusing or misleading to a novice. The "Specialized" part in Figure 2 presents a good example.

Cultural. Some context can be interpreted differently depending on the cultural background of the interpreter, as words, symbols, and actions with specific meanings and connotations. Cultural norms, values, and historical experiences shape how individuals understand and react to information. This can lead to significant textual ambiguities where meaning is lost or distorted. The "Cultural" part in Figure 2 offers a good example.

Dual-ambiguity. This data type highlights our unique contribution. In this type, the context and the image are both ambiguous, but their combined information allows for clear disambiguation. Figure 2 gives a good example. The context ambiguity covers the six former types. Regarding image ambiguity, following CODIS (Luo et al., 2024), image ambiguity can be further categorized into distinct types, such as location and orientation, temporal information, cultural background, attributes, and relationships. However, for the purpose of this paper, we group these various types under the general term "image ambiguity". Thus, dual-ambiguity specifically denotes the situation where both the context and the image exhibit ambiguity.

333

334

335

336

337

338

339

340

341

342

300

301

302

303

304

305



Figure 3: Distribution of seven categories (left) and scenarios (right) of our benchmark.

3.2 Instruction Design

263

265

269

270

271

272

273

274

277

278

284

287

290

291

292

In order to ensure that model fully understands the context and image instead of making choices randomly, we organize our dataset in pairs. For the first six data types, the query can be represented as $(\mathcal{C}, \mathcal{Q}, \mathcal{I}_i)$. Each pair consists of an identical ambiguous context C and a question Q, which are presented alongside i different unambiguous images $(\mathcal{I}_1, \mathcal{I}_2, \ldots, \mathcal{I}_i)$.

For the dual-ambiguity data type, we manually group queries that look similar. Even within a pair, the context C_i , the question Q_i , and the image I_i may differ for each query instance. More formally, each pair can be represented as (C_i, Q_i, I_i) .

We give MLLMs these *i* queries separately and get the model outputs $(\mathcal{O}_1, \mathcal{O}_2, \dots, \mathcal{O}_i)$. These outputs are subsequently compared with the ground truth answers (A_1, A_2, \ldots, A_i) .

Data Collection 3.3

In this section, we introduce our process for constructing this dataset. Our data collection process includes three steps.

Context Collection. We manually collected ambiguous contexts which can only be resolved with external images. Most of the contexts come from the Internet, while others are created manually. These contexts span three languages: Chinese, English, and Malay. We meticulously reviewed every context to exclude any that were unambiguous. In total, we collected 1278 contexts. The first six types comprise 906 contexts, and the remaining 372 contexts belong to the dual-ambiguity type.

Design of Questions, Images and Answers. For each context, we manually wrote questions and answers. As for images, we also manually col-296 lected them from the Internet or designed them our-297 selves. Specifically, for the dual-ambiguity type, the majority of images are sourced from CODIS 299

(Luo et al., 2024), a benchmark constructed using ambiguous images. The data are compiled with following rules:

(1) Questions are designed to target ambiguous aspects within the contexts. Disambiguation of these contexts is not possible without the inclusion of external images.

(2) For each context and question, every unique image associated with them should lead to a distinct interpretation of the context, resulting in a unique answer per image. Crucially, the answer cannot be determined from the image or the context in isolation. The answer can be determined only when we give MLLMs the query in the format $(\mathcal{C}, \mathcal{Q}, \mathcal{I})$.

(3) To balance the performance and evaluation efficiency, MLLMs were required to generate outputs following a specific template with clear, predefined options. This method ensures the objectivity of responses and facilitates their efficient evaluation. Data Verification. Five annotators participated in this process. To ensure dataset quality, each submission was cross-checked by the remaining four annotators. Data were retained only if they satisfied the following conditions: (1) correctness, (2) distinctiveness from existing data, and (3) compliance with all predefined criteria. Submissions that did not meet these conditions were returned to the annotator for revision.

Finally, our benchmark comprises 1278 queries and 501 ($\mathcal{C}, \mathcal{Q}, \mathcal{I}$) pairs, categorized into seven types. The distribution of categories and scenarios is illustrated in Figure 3.

Evaluation Method 3.4

For the k-th pair of queries, we decide to use $(\mathcal{O}_{k1}, \mathcal{O}_{k2}, \ldots, \mathcal{O}_{ki})$ to denote model outputs of a pair of queries, and $(A_{k1}, A_{k2}, \ldots, A_{ki})$ to denote the groundtruth answers. The evaluation of model outputs can be represented as follows:

$$\operatorname{Eval}(\mathcal{O}_{ki}) = \begin{cases} 1 & \text{if } \mathcal{O}_{ki} \text{ matches } \mathcal{A}_{ki} \\ 0 & \text{otherwise} \end{cases}, i \in \mathbb{Z}^+$$

Following Fu et al. (2023a), we leverage two metrics, pair-wise accuracy Acc_p and query-wise accuracy Acc_q for our evaluation metrics, which are calculated as follows:

 n_k

$$\operatorname{Acc}_{p} = \frac{1}{n_{p}} \sum_{k=1}^{n_{p}} \prod_{i=1}^{n_{k}} \operatorname{Eval}(\mathcal{O}_{ki}),$$
343

$$\operatorname{Acc}_{q} = \frac{1}{n_{q}} \sum_{k=1}^{n_{p}} \sum_{i=1}^{n_{k}} \operatorname{Eval}(\mathcal{O}_{ki}).$$
345

(a) Without confusing options																
Model	Poly Acc _p	semy Acc _q	Homo Acc _p	onymy Acc _q	Grar Acc _p	nmar Acc _q	Sem: Acc _p	antics Acc _q	Speci Acc _p	alized Acc _q	Cult Acc _p	tural Acc _q	Dual-a Acc _p	mbiguity Acc _q	Ove Acc _p	erall Acc _q
						API-b	ased Mo	dels								
GPT-4V	35.23	55.87	32.48	53.20	43.93	65.89	37.21	60.71	29.55	60.00	19.05	38.82	20.62	42.74	30.39	52.70
GPT-40	36.36	61.45	44.44	64.00	41.12	65.42	44.19	64.29	47.73	70.53	28.57	47.06	21.65	53.23	34.96	60.13
Gemini-2.0-flash	19.48	52.87	38.18	63.14	17.78	50.56	19.51	50.00	30.77	63.53	25.00	46.91	12.36	36.16	23.46	52.11
Claude-5.5-Solillet	5.41	10.70	2.50	10.00	5.01	21.05	0.00	14.27	0.82	17.09	0.00	5.88	23.20	51.10	9.51	20.30
		26.21	5.10	46.40	0	pen-sour		eis > 7 B	0.07	26.04	0.00	47.06	14.00	52.07	1 4 7 4	45.15
Kimi-VL Llama-3.2-Vision-11B	3.37	30.31 63.60	5.13 37.61	46.40 66 80	25.03	43.46	2.38	40.48 57.14	2.27	30.84 63.16	0.00	47.06	14.08	52.27	4.74	45.15
MiniCPM-o 2.6	38.20	66.48	38.46	66.40	33.33	63.55	28.57	60.71	40.91	69.47	55.29	29.27	4.23	20.53	34.00	53.10
Idefics3-8B-Llama3	39.33	67.04	37.31	62.80	31.48	62.62	35.71	65.48	31.32	56.84	24.09	56.47	19.72	49.87	32.38	58.91
InternVL2-8B	20.24	49.56	23.76	53.67	10.29	40.44	23.81	47.62	20.00	47.13	21.95	48.24	20.97	47.58	20.08	48.16
InternVL2.5-8B-MPO	42.26	69.32	43.56	68.81	25.00	60.29	42.86	69.05	37.50	58.62	19.51	57.65	24.19	51.88	35.73	61.68
InternVL2.5-8B-MPO-AWQ	42.86	69.91	46.53	70.18	30.88	64.71	33.33	61.90	37.50	60.92	21.95	60.00	25.81	54.30	37.32	63.32
Open-source Models ≤ 7B																
Deepseek-VL-Tiny	0.00	44.69	5.98	51.20	0.00	48.60	2.33	48.81	6.82	47.37	0.00	49.41	16.49	54.03	6.77	50.12
LLaVA-v1.6-vicuna-7b	15.73	54.19	21.37	54.80	14.81	54.21	21.43	55.95	11.36	42.11	12.20	51.76	16.90	50.40	16.46	52.26
Owen2 5-VI -3B-Instruct	33.03	55.87 64.60	23.08	56.88	12.90	53.68	10.07	58.55 50.00	22 50	49.47	2.44	47.00	20.90	30.40 40.86	25 55	53.60 53.56
Owen2.5-VL-7B-Instruct	34.52	60.47	32.67	60.09	29.41	58.09	33.33	50.00	35.00	63.22	14.63	54.12	19.35	49.19	29.94	56.29
mPLUG-Owl3-7B-240728	32.74	63.42	35.64	63.30	22.06	58.82	28.57	61.90	22.50	55.17	48.24	14.63	27.42	52.96	28.74	58.25
mPLUG-Owl3-2B-241014	33.73	65.68	36.08	64.71	27.47	62.09	30.95	62.20	30.00	60.92	10.00	55.74	16.49	51.61	25.65	59.29
LLaVA-OneVision	30.33	61.45	38.46	64.80	30.56	61.22	28.57	58.33	27.27	55.79	21.95	60.00	9.86	47.20	28.37	57.20
(b) With confusing options																
Model	Poly	semy Acc.	Home Acc.	onymy Acc.	Grar Acc.	nmar Acc.	Sema Acc.	antics Acc.	Speci Acc.	alized Acc.	Cult Acc.	tural Acc.	Dual-a Acc.	nbiguity Acc.	Ove	erall Acc.
	- · · P				· · P	API-b	ased Mo	dels	- · · P		· · · p		P	4	- · · P	
CPT AV	26.51	47.03	20.62	41.18	26.27	52.20	21 /2	45.12	20.00	47.13	6.67	26.23	15.46	30.65	10.03	40.45
GPT-40	32.53	50.30	25.77	45.10	31.87	56.04	23.81	46.34	20.00	51.72	13.33	20.23	17.53	36.29	23.92	40.43
Gemini-2.0-flash	8.06	36.80	8.82	32.41	14.06	44.53	8.70	40.91	4.76	25.58	6.67	32.26	11.02	31.23	10.00	34.85
Claude-3.5-Sonnet	3.61	22.49	4.12	22.55	4.40	23.08	2.38	20.73	2.50	27.59	3.33	13.11	15.59	35.64	7.56	26.50
					0	pen-sour	ce Mode	ls > 7B								
Kimi-VL	1.13	31.84	1.71	39.60	1.85	40.19	0.00	36.90	0.00	32.63	0.00	41.18	5.71	18.82	1.50	31.98
Llama-3.2-Vision-11B	21.35	56.98	37.61	63.60	13.89	43.93	14.29	45.24	11.37	48.42	17.07	54.12	5.71	34.67	19.87	48.00
MiniCPM-o 2.6	20.22	53.07	34.18	60.40	28.70	58.41	30.95	59.52	18.18	60.00	14.63	50.58	2.86	12.63	23.50	44.40
Idefics3-8B-Llama3	29.21	58.10 35.20	39.32 8.55	62.80	24.07	55.14 34.58	35.71	60.71 36.00	31.82	62.11	7 32	51.76	12.86	50.54 31.45	28.14	34.24
					/.+1	54.56	/.14	50.90	20.45	41.05	1.32	50.47	8.00	51.45	9.70	52.01
InternVL2 5-8B-MPO	26.97	58 10	27 35	54.80	27.78	56 54	28 57	58 33	25.00	52.63	19.51	5176	20.97	46 51	25.95	1111
InternVL2.5-8B-MPO InternVL2.5-8B-MPO-AWQ	26.97 31.46	58.10 63.69	27.35 33.33	54.80 59.20	27.78 34.26	56.54 62.62	28.57 30.95	58.33 59.52	25.00 31.81	52.63 57.89	19.51 14.63	51.76 52.94	20.97 22.58	46.51 50.27	25.95 29.94	55.01 57.31
InternVL2.5-8B-MPO InternVL2.5-8B-MPO-AWQ	26.97 31.46	58.10 63.69	27.35 33.33	54.80 59.20	27.78 34.26	56.54 62.62 pen-sour	28.57 30.95	58.33 59.52 els $\leq 7B$	25.00 31.81	52.63 57.89	19.51 14.63	51.76 52.94	20.97 22.58	46.51 50.27	25.95 29.94	57.31
InternVL2-58B-MPO InternVL2.5-8B-MPO-AWQ Deepseek-VL-Tiny	26.97 31.46	58.10 63.69 21.30	0.00 27.35 33.33	54.80 59.20 21.08	27.78 34.26 0.00	56.54 62.62 pen-sour 30.77	28.57 30.95 rce Mode 0.00	58.33 59.52 els $\leq 7B$ 24.39	25.00 31.81 5.00	52.63 57.89 31.03	19.51 14.63	51.76 52.94 27.87	20.97 22.58 9.28	46.51 50.27 35.48	25.95 29.94	57.31 28.61
Intern VL2-5-BB-MPO Intern VL2.5-8B-MPO-AWQ Deepseek-VL-Tiny LLaVA-v1.6-vicuna-7b	26.97 31.46	58.10 63.69 21.30 51.95	0.00 27.35 33.33 0.00 19.66	54.80 59.20 21.08 58.00	27.78 34.26 0.00 14.81	56.54 62.62 pen-sour 30.77 56.07	28.57 30.95 ce Mode 0.00 19.05	$58.33 \\ 59.52 \\ els \le 7B \\ 24.39 \\ 54.76 \\ els \le 7B \\ 24.39 \\ 54.76 \\ els \le 7B \\ 24.39 \\ 54.76 \\ els \le 7B \\ $	25.00 31.81 5.00 11.36	52.63 57.89 31.03 49.47	19.51 14.63 0.00 14.63	27.87 55.29	20.97 22.58 9.28 11.43	46.51 50.27 35.48 46.77	25.95 29.94 3.47 15.19	28.61 52.54
Intern VL2-5-BB-MPO Intern VL2.5-8B-MPO Intern VL2.5-8B-MPO-AWQ Deepseek-VL-Tiny LLaVA-v1.6-vicuna-7b LLaVA-v1.6-mistral-7b	13.48 26.97 31.46 0.00 13.48 19.10	58.10 63.69 21.30 51.95 57.54	0.00 27.35 33.33 0.00 19.66 23.93	54.80 59.20 21.08 58.00 59.20	27.78 34.26 0 0.00 14.81 12.96	56.54 62.62 pen-sour 30.77 56.07 54.21	28.57 30.95 cce Mode 0.00 19.05 26.19	$58.33 \\ 59.52$ $els \le 7B$ $24.39 \\ 54.76 \\ 69.91$	25.00 31.81 5.00 11.36 18.18	52.63 57.89 31.03 49.47 51.58	19.51 14.63 0.00 14.63 2.44	27.87 55.29 49.41	20.97 22.58 9.28 11.43 14.29	46.51 50.27 35.48 46.77 47.58	25.95 29.94 3.47 15.19 17.14	28.61 52.54 53.71
Intern VL2.5-8B-MPO Intern VL2.5-8B-MPO Intern VL2.5-8B-MPO-AWQ Deepseek-VL-Tiny LLaVA-v1.6-vicuna-7b LLaVA-v1.6-mistral-7b Qwen2.5-VL-3B-Instruct	0.00 13.48 19.10 25.84	58.10 58.10 63.69 21.30 51.95 57.54 56.42	0.00 27.35 33.33 0.00 19.66 23.93 28.21	54.80 59.20 21.08 58.00 59.20 58.00 59.20	27.78 34.26 0 0.00 14.81 12.96 25.00 27.76	56.54 62.62 pen-sour 30.77 56.07 54.21 57.94	28.57 30.95 ce Mode 0.00 19.05 26.19 19.05	58.33 59.52 24.39 54.76 69.91 53.57 51.10	25.00 31.81 5.00 11.36 18.18 11.36 27.27	52.63 57.89 31.03 49.47 51.58 50.53 (2.00)	19.51 14.63 0.00 14.63 2.44 14.63	27.87 55.29 49.41 42.35	20.97 22.58 9.28 11.43 14.29 24.19 20.07	46.51 50.27 35.48 46.77 47.58 45.43	25.95 29.94 3.47 15.19 17.14 22.95 20.76	28.61 52.54 52.23
Intern VL2-5-BB-MPO Intern VL2.5-8B-MPO Intern VL2.5-8B-MPO-AWQ Deepseek-VL-Tiny LLaVA-v1.6-vicuna-7b LLaVA-v1.6-wistral-7b Qwen2.5-VL-7B-Instruct Qwen2.5-VL-7B-Instruct Qwen2.5-VL-7B-Instruct	26.97 31.46 0.00 13.48 19.10 25.84 13.48 21.35	53.10 58.10 63.69 21.30 51.95 57.54 56.42 42.46 56.42	0.00 27.35 33.33 0.00 19.66 23.93 28.21 19.66 30.77	21.08 58.00 59.20 21.08 58.00 59.20 58.00 50.80 60.00	27.78 34.26 0 0.00 14.81 12.96 25.00 27.78 22.22	56.54 62.62 pen-sour 30.77 56.07 54.21 57.94 55.61 56.54	28.57 30.95 rce Mode 0.00 19.05 26.19 19.05 26.19 26.19 26.19	$58.33 \\ 59.52 \\ \hline 24.39 \\ 54.76 \\ 69.91 \\ 53.57 \\ 51.19 \\ 59.52 \\ \hline$	25.00 31.81 5.00 11.36 18.18 11.36 27.27 18.18	52.63 57.89 31.03 49.47 51.58 50.53 60.00 53.68	19.51 14.63 0.00 14.63 2.44 14.63 9.76 9.76	27.87 55.29 49.41 42.35 44.71 43.53	20.97 22.58 9.28 11.43 14.29 24.19 20.97 24.19	46.51 50.27 35.48 46.77 47.58 45.43 48.66 47.85	25.95 29.94 3.47 15.19 17.14 22.95 20.76 23.35	28.61 52.54 53.71 52.23 50.12 53.79
Deepseek-VL-Tiny LLaVA-v1.6-vicuna-7b LLaVA-v1.6-vicuna-7b LLaVA-v1.6-mistral-7b Qwen2.5-VL-7B-Instruct Qwen2.5-VL-7B-Instruct mPLUG-0wl3-2B-241014	13.43 26.97 31.46 0.00 13.48 19.10 25.84 13.48 21.35 24.10	58.10 58.10 63.69 21.30 51.95 57.54 56.42 42.46 56.42 53.85	0.00 27.35 33.33 0.00 19.66 23.93 28.21 19.66 30.77 30.93	21.08 58.00 59.20 21.08 58.00 59.20 58.00 50.80 60.00 60.78	27.78 34.26 0 0.00 14.81 12.96 25.00 27.78 22.22 25.27	56.54 62.62 pen-sour 30.77 56.07 54.21 57.94 55.61 56.54 59.34	28.57 30.95 ce Mode 0.00 19.05 26.19 19.05 26.19 26.19 26.19 33.33	$58.33 \\ 59.52 \\ \hline els \le 7B \\ 24.39 \\ 54.76 \\ 69.91 \\ 53.57 \\ 51.19 \\ 59.52 \\ 58.54 \\ \hline$	25.00 31.81 5.00 11.36 18.18 11.36 27.27 18.18 32.50	52.63 57.89 31.03 49.47 51.58 50.53 60.00 53.68 60.92	19.51 14.63 0.00 14.63 2.44 14.63 9.76 9.76 16.67	27.87 55.29 49.41 42.35 44.71 43.53 49.18	20.97 22.58 9.28 11.43 14.29 24.19 20.97 24.19 13.40	46.51 50.27 35.48 46.77 47.58 45.43 48.66 47.85 46.77	25.95 29.94 3.47 15.19 17.14 22.95 20.76 23.35 22.70	28.61 52.54 53.71 52.23 50.12 53.79 54.28
Intern VL2-5-BB-MPO Intern VL2.5-8B-MPO Intern VL2.5-8B-MPO-AWQ Deepseek-VL-Tiny ILaVA-v1.6-vicuna-7b ILaVA-v1.6-wistral-7b Qwen2.5-VL-3B-Instruct Qwen2.5-VL-7B-Instruct Qwen2.5-VL-7B-Instruct mPLUG-0W13-2B-241014 ILaVA-One Vision	13.48 26.97 31.46 0.00 13.48 19.10 25.84 13.48 21.35 24.10 40.44 10.44	21.30 57.54 56.42 42.46 56.42 53.85 68.16	27.35 33.33 0.00 19.66 23.93 28.21 19.66 30.77 30.93 39.32	21.08 59.20 21.08 58.00 59.20 58.00 50.80 60.00 60.78 65.60	27.78 34.26 0.00 14.81 12.96 25.00 27.78 22.22 25.27 34.26	56.54 62.62 pen-sour 30.77 56.07 54.21 57.94 55.61 56.54 59.34 63.08	28.57 30.95 rce Mode 0.00 19.05 26.19 19.05 26.19 26.19 33.33 33.33	$58.33 \\ 59.52 \\ \hline \\ $	25.00 31.81 5.00 11.36 18.18 11.36 27.27 18.18 32.50 31.82	52.63 57.89 31.03 49.47 51.58 50.53 60.00 53.68 60.92 60.00	19.51 14.63 0.00 14.63 2.44 14.63 9.76 9.76 16.67 17.07	27.87 55.29 49.41 42.35 44.71 43.53 49.18 55.29	20.97 22.58 9.28 11.43 14.29 24.19 20.97 24.19 13.40 8.57	46.51 50.27 35.48 46.77 47.58 45.43 48.66 47.85 46.77 41.67	25.95 29.94 3.47 15.19 17.14 22.95 20.76 23.35 22.70 31.40	28.61 52.54 53.71 52.23 50.12 53.79 54.28 57.39

Table 2: Results of MLLMs on MUCAR benchmark under two settings: (top) without confusing options and (bottom) with confusing options.

where n_k is the number of queries in each pair, n_p is number of data pairs, n_q is the total number of individual queries. For Acc_p, models score only if their answers to a pair of queries are both correct. For Acc_q, models score for each single query they answer correctly.

4 Experiments

4.1 Models

347

357

360

We evaluate a total of 19 models covering a range of scales and architectures. Our evaluated proprietary models include GPT-4V (OpenAI, 2023), GPT-40 (OpenAI, 2024), Gemini (Gemini Team et al., 2023), and Claude-3.5-Sonnet (Anthropic, 2024). For open-source models, we include Deepseek-VL-Tiny (Lu et al., 2024), Kimi-VL (Team et al., 2025),

Llama-3.2-Vision-11B (Meta, 2024), MiniCPM-o 2.6 (Yao et al., 2024), InternVL2.5 series (Wang et al., 2024b; Chen et al., 2024c), LLaVA-v1.6vicuna-7b (Liu et al., 2024), Qwen2.5-VL series (Team, 2025). Details of these models are listed in Table 7 in Appendix B.

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

4.2 Main Results

Main experimental results on our benchmark of all 19 models are reported in Table 2.

Overall Performance. Across all evaluated models, InternVL2.5-8B-MPO-AWQ achieves the best overall accuracy ($Acc_q = 67.03\%$), followed closely by MiniCPM-0 2.6 (66.48%) and InternVL2.5-8B-MPO (65.71%). Among proprietary models, GPT-40 outperforms the others, obtaining an overall accuracy of 53.91%, slightly higher than GPT-



Figure 4: Accuracy under different input settings. C = Context, Q = Question, I = Image. For example, (C, Q, I) includes all three; (C, Q) and (C, I) include only the specified components.

4V (52.70%). In contrast, Claude-3.5-Sonnet and Kimi-VL underperform, showing limited ability in disambiguation tasks.

Results with Different Model Size. Open-source models with scales *larger than 7B* generally outperform smaller ones, with all top-performing models falling within this range of scale, which is likely to benefit from richer training data and more advancedadvanced architectures. In comparison, models *smaller than or equal to 7B* show a clear performance gap. Although certain models, such as LaVA-Onevision (61.30%) and mPLUG-Owl2-3B (58.21%), perform competitively, most smaller models struggle with complex ambiguities, particularly in semantic and cultural contexts.

380

382

386

390

397

400

401

402

403

404

405

406

407

Results on Different Categories. We further break down the results by disambiguation categories, and find that InternVL2.5-8B-MPO-AWQ consistently leads in most categories, especially in Homonymy (70.18%), Grammar (64.71%), and Semantics (61.90%). Notably, MiniCPM-0 2.6 excels in the Specialized category (69.47%), suggesting domain knowledge plays a key role. In the Cultural category, which requires understanding crosscultural references, models like GPT-40 (34.12%) and InternVL2.5-8B-MPO-AWQ (60.00%) show relatively stronger performance. On the other hand, most models perform poorly in the Polysemy and Dual-Ambiguity categories, reflecting the inherent challenges in resolving subtle or cross-modal ambiguities.

In summary, model size and architecture signif-408 icantly affect cross-modal disambiguation perfor-409 mance. Larger models and instruction-finetuned 410 models usually demonstrate better generalization. 411 The disparity across categories reveals the diverse 412 challenges in context-dependent reasoning, partic-413 ularly in categories involving semantic, cultural, or 414 compound ambiguities. 415

Madal	Ov	erall
Wodel	Acc _p	Acc_q
Qwen2.5-VL-7B-Instruct (w/)	20.73	50.72
Qwen2.5-VL-7B-Instruct (w/o)	7.06 (13.67)	12.57 (138.15)
InternVL2.5-8B-MPO-AWQ (w/)	30.98	60.20
InternVL2.5-8B-MPO-AWQ (w/o)	4.56 (126.42)	8.38 (151.82)
LLaVA-OneVision (w/)	34.62	63.84
LLaVA-OneVision (w/o)	0.00 (134.62)	0.66 (163.18)
MiniCPM-o 2.6 (w/)	26.42	57.44
MiniCPM-o 2.6 (w/o)	1.14 (↓25.28)	4.63 (152.81)

Table 3: Ablation study: Only input (Q, I), with confusion options. Performance drop (\downarrow) indicates the gap compared to full input.

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

5 Analyses and Discussion

5.1 Ablation on Benchmark

Figure 4 reports the accuracy of four models under three input settings: (C, Q, I), (C, I), and (C, Q). All models achieve the highest accuracy with full input, with LLaVA-OneVision reaching 63.84%, followed by InternVL2.5 (60.20%), MiniCPM-o 2.6 (57.44%), and Qwen2.5 (50.72%). Removing the question while keeping context and image (C, I) results in moderate drops—for example, LLaVA drops to 55.58%, and InternVL2.5 to 21.41%. In contrast, removing the image (C, Q)leads to drastic degradation: LLaVA drops to 0.00%, and Qwen2.5 to 3.42%. This highlights the essential role of image information in resolving multimodal ambiguity.

Table 3 presents the ablation results using only question and image inputs, with confusing options included. All models exhibit significant performance degradation when context is removed. For example, LLaVA-OneVision drops from 63.84% to 0.66% in Acc_q (\downarrow 63.18), and from 34.62% to 0.00% in Acc_p (\downarrow 34.62). Similar trends are observed for InternVL2.5 and MiniCPM-o, which also suffer large drops in both metrics. These results underscore the importance of contextual information in resolving ambiguity, especially in the presence of visually or semantically confusing alternatives.

5.2 Discussion

As shown in Figure 5, this example illustrates how the interpretation of the phrase "666" is highly dependent on cultural and visual context, highlighting the necessity of cross-modal disambiguation. **Scenario 1:** The accompanying image shows the Forbidden City in Beijing, indicating a modern Chinese cultural context. In this setting, "666" is widely used as internet slang to express praise, meaning "awesome" or "skillful." *Answer: Positive.* **Scenario 2:** The image depicts a European



Figure 5: Pipeline of our proposed framework.

Gothic cathedral—Notre-Dame de Paris—evoking a Western Christian context. Here, "666" is traditionally associated with the "number of the beast" from the Bible, conveying a negative connotation. *Answer: Negative*.

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

480

481

482

483

484

5.3 Further Exploration: An Agent Framework for Ambiguity Resolution

To better address the cross-modal ambiguity resolution, we propose an agent-based framework, as shown in Figure 5. First, the model takes the task description, question, and ambiguous context/image as input to identify the ambiguity and its significance. Second, it generates contextual evidence to resolve the ambiguity. Finally, the model bridges the context and image through logical reasoning to produce the final answer.

Formally, given a task description \mathcal{D} , a question \mathcal{Q} , and a multimodal ambiguous context \mathcal{X} (e.g., an image and textual), the agent-based model \mathcal{M} solves the task through a three-step process as illustrated in Figure 5:

Step 1: Ambiguity Detection. The model first detects the ambiguity and explains its significance:

$$\mathcal{A} = \mathcal{M}(\mathcal{D} \oplus \mathcal{Q} \oplus \mathcal{X} \oplus \mathcal{T}_1)$$

where \mathcal{T}_1 is a prompt guiding the model to identify potential ambiguity and why it matters, and \mathcal{A} denotes the ambiguity explanation.

Step 2: Contextual Evidence Extraction. The model then extracts relevant contextual evidence to

Madal	CO	DIS	MMA	MU	CAR	Ove	rall
Wiodei	Accp	Accq	Acc	Acc _p	Accq	Accp	Accq
Vanilla	36.26	59.49	72.0	32.35	53.91	34.31	54.61
CoT	36.81	60.76	68.0	28.57	56.41	32.69	52.44
CODIS	36.80	60.30	71.0	28.29	56.96	32.55	53.19
OURS	42.49	63.46	84.0	44.87	66.78	43.68	64.11

Table 4: Results of our method on CODIS, MMA and MUCAR based on GPT-40-2024-11-20.

resolve the ambiguity:

$$\mathcal{E} = \mathcal{M}(\mathcal{D} \oplus \mathcal{Q} \oplus \mathcal{X} \oplus \mathcal{A} \oplus \mathcal{T}_2),$$

485

487

488

489

490

491

492

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

where T_2 instructs the model to generate explanatory evidence from the context, and \mathcal{E} denotes the extracted evidence.

Step 3: Logical Reasoning and Final Answer. Finally, the model uses evidence to logically align context and image, and generate the final answer:

$$\mathcal{R} = \mathcal{M}(\mathcal{D} \oplus \mathcal{Q} \oplus \mathcal{X} \oplus \mathcal{E} \oplus \mathcal{T}_3),$$
 493

where T_3 prompts the model to conduct reasoning and provide the answer \mathcal{R} .

Unified Prompt Alternatively, the entire process can be completed with a unified prompt:

$$\mathcal{A}, \mathcal{E}, \mathcal{R} = \mathcal{M}(\mathcal{D} \oplus \mathcal{Q} \oplus \mathcal{X} \oplus \mathcal{T}),$$

where $\mathcal{T} = \mathcal{T}_1 \oplus \mathcal{T}_2 \oplus \mathcal{T}_3$.

To demonstrate the generality of our framework across different types of ambiguity, we also report performance of our framework on other wide discussed benchmarks, including CODIS (Luo et al., 2024) and MMA (Wang et al., 2024a). As shown in Table 4, our method consistently outperforms all baselines across benchmarks. Compared to CoT and CODIS-specific prompting, our approach is not only simple but also effective, demonstrating strong potential across different types of ambiguity.

6 Conclusion

We present MUCAR, a benchmark designed to evaluate MLLMs in resolving ambiguities across visual, textual, and multilingual contexts. Unlike prior benchmarks, MUCAR targets crossmodal disambiguation through multilingual text and image-text ambiguity cases. Evaluation of 19 state-of-the-art MLLMs reveals a clear gap from human-level performance, highlighting the need for more context-aware and cross-modally grounded models. We also introduce a simple agent-based framework that improves disambiguation through explicit reasoning. MUCAR aims to guide future research toward more robust and interpretable multimodal systems in this direction.

Limitations

525

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

556

557

558

559

561

562

563

564

566

567

568

571

572

574

While MUCAR offers a novel and rigorous benchmark for multilingual cross-modal ambiguity res-527 olution, it has several limitations. It covers only 528 three languages, limiting generalizability to low-529 resource or typologically diverse languages. The 530 531 curated examples may not capture the complexity and noise of real-world multimodal data. Its partial 532 reliance on GPT-based evaluation introduces potential biases, and the agent-based reasoning framework, though effective in structured tasks, may 535 struggle with open-ended scenarios. Future work 536 should address broader linguistic coverage, realworld settings, and more robust reasoning methods.

References

- Anthropic. 2024. Claude 3: A new generation of ai. Accessed on May 5, 2025.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. VQA: Visual question answering. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2425–2433.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, and Feng Zhao. 2024a. Are we on the right way for evaluating large vision-language models?
- Xiang Chen, Chenxi Wang, Yida Xue, Ningyu Zhang, Xiaoyan Yang, Qiang Li, Yue Shen, Lei Liang, Jinjie Gu, and Huajun Chen. 2024b. Unified hallucination detection for multimodal large language models.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. 2015. Microsoft COCO Captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. 2024c. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. 2024d. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*.
- Zhenfang Chen, Qinhong Zhou, Yikang Shen, Yining Hong, Hao Zhang, and Chuang Gan. 2023. See, think, confirm: Interactive prompting between vision

and language models for knowledge-based visual reasoning. *arXiv preprint arXiv:2301.05226*.

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

- Jiwan Chung, Seungwon Lim, Jaehyun Jeon, Seungbeen Lee, and Youngjae Yu. 2024. Can visual language models resolve textual ambiguity with visual cues? let visual puns tell you! In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 2452–2469, Miami, Florida, USA. Association for Computational Linguistics.
- Chenhang Cui, Yiyang Zhou, Xinyu Yang, Shirley Wu, Linjun Zhang, James Zou, and Huaxiu Yao. 2023. Holistic analysis of hallucination in GPT-4V(ision): Bias and interference challenges. *arXiv preprint arXiv:2311.03287*.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. InstructBLIP: Towards general-purpose visionlanguage odels with instruction tuning. *arXiv preprint arXiv:2305.06500*.
- Rachel N Denison, William T Adler, Marisa Carrasco, and Wei Ji Ma. 2018. Humans incorporate attentiondependent uncertainty into perceptual decisions and confidence. *Proceedings of the National Academy of Sciences*, 115(43):11090–11095.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. 2023a. MME: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*.
- Chaoyou Fu, Renrui Zhang, Haojia Lin, Zihan Wang, Timin Gao, Yongdong Luo, Yubo Huang, Zhengye Zhang, Longtian Qiu, Gaoxiang Ye, et al. 2023b. A challenger to GPT-4V? early explorations of Gemini in visual expertise. *arXiv preprint arXiv:2312.12436*.
- Lingyu Gao, Aditi Chaudhary, Krishna Srinivasan, Kazuma Hashimoto, Karthik Raman, and Michael Bendersky. 2024. Ambiguity-aware in-context learning with large language models.
- G Gemini Team, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Tianrui Guan, Fuxiao Liu, Xiyang Wu Ruiqi Xian Zongxia Li, Xiaoyu Liu Xijun Wang, Lichang Chen Furong Huang Yaser Yacoob, and Dinesh Manocha Tianyi Zhou. 2023. HallusionBench: An advanced diagnostic suite for entangled language hallucination & visual illusion in large vision-language models. *arXiv e-prints*, pages arXiv–2310.
- Tanmay Gupta and Aniruddha Kembhavi. 2023. Visual Programming: Compositional visual reasoning without training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14953–14962.

- 632 633 637 639 647

- 653

679

680

683

Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. 2024. Building and better understanding vision-language models: insights and future directions.

- Youngwon Lee, Seung won Hwang, Ruofan Wu, Feng Yan, Danmei Xu, Moutasem Akkad, Zhewei Yao, and Yuxiong He. 2025. Agentic verification for ambiguous query disambiguation.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024. Llavaonevision: Easy visual task transfer.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023. SEED-Bench: Benchmarking multimodal llms with generative comprehension. arXiv preprint arXiv:2307.16125.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning. arXiv preprint arXiv:2310.03744.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024. Improved baselines with visual instruction tuning.
- Yang Liu, Guanbin Li, and Liang Lin. 2023b. Crossmodal causal relational reasoning for event-level visual question answering. IEEE Transactions on Pattern Analysis and Machine Intelligence.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2023c. MM-Bench: Is your multi-modal model an all-around player? arXiv preprint arXiv:2307.06281.
- Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan. 2024. Deepseek-vl: Towards real-world vision-language understanding.
- 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. 2024. CODIS: Benchmarking context-dependent visual comprehension for multimodal large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 10639–10659, Bangkok, Thailand. Association for Computational Linguistics.
- Jianjie Luo, Yehao Li, Yingwei Pan, Ting Yao, Jianlin Feng, Hongyang Chao, and Tao Mei. 2023. Semantic-conditional diffusion networks for image captioning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 23359-23368.

Xinyu Ma, Xuebo Liu, Derek F Wong, Jun Rao, Bei Li, Liang Ding, Lidia S Chao, Dacheng Tao, and Min Zhang. 2024. 3AM: An ambiguity-aware multimodal machine translation dataset. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 1-13, Torino, Italia. ELRA and ICCL.

685

686

688

689

690

691

692

693

694

695

696

697

698

699

700

702

703

704

705

706

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

733

734

735

736

738

- Meta. 2024. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models.
- Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. AmbigQA: Answering ambiguous open-domain questions. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 5783– 5797, Online. Association for Computational Linguistics.
- Heejeong Nam, Jinwoo Ahn, Keummin Ka, Jiwan Chung, and Youngjae Yu. 2024. VAGUE: Visual Contexts Clarify Ambiguous Expressions. arXiv eprints, page arXiv:2411.14137.
- Roberto Navigli. 2009. Word sense disambiguation: A survey. ACM Comput. Surv., 41(2).
- OpenAI. 2023. GPT-4V(ision) System Card.
- OpenAI. 2024. Hello GPT-40. Accessed on May 5, 2025.
- Mohammadmostafa Rostamkhani, Baktash Ansari, Hoorieh Sabzevari, Farzan Rahmani, and Sauleh Eetemadi. 2024. Illusory vqa: Benchmarking and enhancing multimodal models on visual illusions.
- Zhenwei Shao, Zhou Yu, Meng Wang, and Jun Yu. 2023. Prompting large language models with answer heuristics for knowledge-based visual question answering. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 14974-14983.
- Weiwei Sun, Hengyi Cai, Hongshen Chen, Pengjie Ren, Zhumin Chen, Maarten de Rijke, and Zhaochun Ren. 2023. Answering ambiguous questions via iterative prompting. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7669-7683, Toronto, Canada. Association for Computational Linguistics.
- Kimi Team, Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chenzhuang Du, Chu Wei, Congcong Wang, Dehao Zhang, Dikang Du, Dongliang Wang, Enming Yuan, Enzhe Lu, Fang Li, Flood Sung, Guangda Wei, Guokun Lai, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haoning Wu, Haotian Yao, Haoyu Lu, Heng Wang, Hongcheng Gao, Huabin Zheng, Jiaming Li, Jianlin Su, Jianzhou Wang, Jiaqi Deng, Jiezhong Qiu, Jin Xie, Jinhong Wang, Jingyuan Liu, Junjie Yan, Kun Ouyang, Liang Chen, Lin Sui, Longhui Yu, Mengfan Dong, Mengnan

Dong, Nuo Xu, Pengyu Cheng, Qizheng Gu, Runjie Zhou, Shaowei Liu, Sihan Cao, Tao Yu, Tianhui Song, Tongtong Bai, Wei Song, Weiran He, Weixiao Huang, Weixin Xu, Xiaokun Yuan, Xingcheng Yao, Xingzhe Wu, Xinxing Zu, Xinyu Zhou, Xinyuan Wang, Y. Charles, Yan Zhong, Yang Li, Yangyang Hu, Yanru Chen, Yejie Wang, Yibo Liu, Yibo Miao, Yidao Qin, Yimin Chen, Yiping Bao, Yiqin Wang, Yongsheng Kang, Yuanxin Liu, Yulun Du, Yuxin Wu, Yuzhi Wang, Yuzi Yan, Zaida Zhou, Zhaowei Li, Zhejun Jiang, Zheng Zhang, Zhilin Yang, Zhiqi Huang, Zihao Huang, Zijia Zhao, and Ziwei Chen. 2025. Kimi-VL technical report.

Qwen Team. 2025. Qwen2.5-vl.

740

741

742

743

744

746

747

748

749

751

753

758

762

771

773

775

778

779

789

790

- Ning Wang, Jiahao Xie, Jihao Wu, Mingbo Jia, and Linlin Li. 2023. Controllable image captioning via prompting. In *Proceedings of the AAAI Conference* on Artificial Intelligence, pages 2617–2625.
- R. Wang, S. Song, L. Ding, S. S. Gu, M. Gong, Y. Iwasawa, and J. Guo. 2024a. Mma: Benchmarking multimodal large language model in ambiguity contexts. *preprint openreview*.
- Weiyun Wang, Zhe Chen, Wenhai Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Jinguo Zhu, Xizhou Zhu, Lewei Lu, Yu Qiao, and Jifeng Dai. 2024b. Enhancing the reasoning ability of multimodal large language models via mixed preference optimization. *arXiv preprint arXiv:2411.10442*.
- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, et al. 2024. Minicpm-v: A gpt-4v level mllm on your phone. arXiv preprint arXiv:2408.01800.
- Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. 2024. mplug-owl3: Towards long image-sequence understanding in multi-modal large language models.
- Kaining Ying, Fanqing Meng, Jin Wang, Zhiqian Li, Han Lin, Yue Yang, Hao Zhang, Wenbo Zhang, Yuqi Lin, Shuo Liu, et al. 2024. Mmt-bench: A comprehensive multimodal benchmark for evaluating large vision-language models towards multitask agi. arXiv preprint arXiv:2404.16006.
- Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6720–6731.

A Details of Datasets

CODIS (Luo et al., 2024) A benchmark to assess the capabilities of MLLMs in context-dependent visual comprehension, consisting of 706 queries over 216 images across five categories. **MMA (Wang et al., 2024a)** A multiple-choice VQA benchmark focused on ambiguity, including 261 questions paired with two divergent images that yield different answers given the same textual context. 793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

MUCAR Table 5 summarizes the distribution of seven ambiguity types across Chinese, English, and Malay. English has the most annotated instances (582), followed by Chinese (364) and Malay (332). *Homonymy* is notably dominant in English, while *Dual-Ambiguity* remains consistently high across all languages. *Cultural* ambiguity appears in Chinese and English but is absent in Malay. *Grammar* and *Polysemy* are relatively balanced, whereas *Specialized* and *Semantics* vary more significantly. These trends reflect both shared and language-specific ambiguity patterns.

Category	Chinese	English	Malay	Total
Polysemy	63	64	52	179
Homonymy	42	176	32	250
Grammar	64	72	78	214
Semantics	24	18	42	84
Specialized	26	60	8	94
Cultural	24	66	0	90
Dual-Amb	126	126	120	372
Overall	364	582	332	1278

 Table 5: Distribution of ambiguity categories across different languages.

Table 6 presents the detailed prompts used during model testing. In the main experiments, we employed English prompts; additionally, we conducted ablation studies using Chinese and Malay prompts to evaluate the impact of different evaluation languages on the experimental results. The table also separately shows the prompts used in the three ablation settings: (Q, I), (C, Q), and C.

B Evaluated Models

We evaluate a total of 19 models covering a range of scales and architectures. Our evaluated proprietary models include GPT-4V (OpenAI, 2023), GPT-4o (OpenAI, 2024), Gemini (Gemini Team et al., 2023), and Claude-3.5-Sonnet (Anthropic, 2024). For open-source models, we include Deepseek-VL-Tiny (Lu et al., 2024), Kimi-VL (Team et al., 2025), Llama-3.2-Vision-11B (Meta, 2024), MiniCPMo 2.6 (Yao et al., 2024), InternVL2-8B (Chen et al., 2024d), InternVL2.5-8B-MPO (Wang et al., 2024b), InternVL2.5-8B-MPO-AWQ (Chen et al., 2024c), LLaVA-v1.6-vicuna-7b (Liu et al., 2024),

	Prompt for Model Inference
Main experiment In English	I'll give you an image. Please answer my question based on the image. Directly select the correct option (A, B, C, D, or E). Use the following format to answer: Answer: [ONLY the option letter; not a complete sentence] Only give me the reply according to this format, don't give me any other words. Now, please answer this question. Question: [QUESTION HERE] Options: [OPTIONS HERE]
Ablation Study In Chinese	我会给你一张图片。请根据图片回答我的问题。直接选择正确选项(A, B, C, D, 或E)。请使用以下格式回答: Answer: [仅为选项字母;不是完整句子] 请只按照此格式回复我,不要给出任何其他文字。现在,请回答这个问题。 问题: [QUESTION HERE] 选项: [OPTIONS HERE]
Ablation Study In Malay	Saya akan berikan anda imej. Sila jawab soalan saya berdasarkan imej tersebut. Pilih terus pilihan yang betul (A, B, C, D, atau E). Gunakan format berikut untuk menjawab: Answer: [HANYA huruf pilihan; bukan ayat penuh] Berikan saya jawapan mengikut format ini sahaja, jangan berikan perkataan lain. Sekarang, sila jawab soalan ini. Soalan: [QUESTION HERE] Pilihan: [OPTIONS HERE]
Ablation Study Only input $(\mathcal{Q}, \mathcal{I})$	I'll give you an image. Please answer my question based on the image. Directly select the correct option (A, B, C, D, or E). Use the following format to answer: Answer: [ONLY the option letter; not a complete sentence] Only give me the reply according to this format, don't give me any other words. Now, please answer this question. Question: [QUESTION HERE] Options: [OPTIONS HERE]
Ablation Study Only input $(\mathcal{C}, \mathcal{Q})$	Please answer my question. Directly select the correct option (A, B, C, D, or E). Use the following format to answer: Answer: [ONLY the option letter; not a complete sentence] Only give me the reply according to this format, don't give me any other words. Now, please answer this question. Question: [QUESTION HERE] Options: [OPTIONS HERE]
Ablation Study Is C ambiguous?	Please determine whether this sentence is ambiguous. If the sentence is ambiguous, please answer 'Yes.'; otherwise, answer 'No.'. Please respond directly with 'yes' or 'no', without any additional content. Sentence: [SENTENCE HERE]

Table 6: Prompt for model inference.

831 Qwen2.5-VL-3B-Instruct (Team, 2025), Qwen2.5-VL-7B-Instruct (Team, 2025). Details of these 832 models are listed in Table 7. Table 7 presents a com-833 prehensive overview of the Multimodal Large Lan-834 guage Models (MLLMs) evaluated in our bench-835 mark. The models are categorized into two groups: API-based models and open-source models. For 837 each model, we list its parameter size category 838 (greater than or less than 7 billion), the vision encoder architecture, the underlying language model 840 (LLM) backbone, and the employed vision-tolanguage (V2L) adapter. API-based models such 842 as GPT-4V and Gemini do not publicly disclose ar-844 chitectural details, while open-source models span a variety of encoders (e.g., SigLIP, CLIP ViT, InternViT), LLM backbones (e.g., Llama, Qwen, In-846 ternLM), and adapter types (e.g., MLP, Linear, 847

XAttn LLM). This table highlights the diversity in architectural design choices across MLLMs.

Table 8 presents the detailed experimental results of multiple models with the language set to Chinese, conducted without any interference options enabled.

Table 9 shows the detailed experimental results of multiple models evaluated with the language set to Chinese, conducted with interference options enabled.

Table 10 provides the detailed experimental results of multiple models evaluated with the language set to Malay, conducted without any interference options enabled.

Table 11 presents the detailed experimental results of multiple models evaluated with the language set to Malay, conducted with interference

Model	Parameters	Vision Encoder	LLM Backbone	V2L Adapter
GPT-4V (OpenAI, 2023)		-	-	-
GPT-40 (OpenAI, 2024)		-	-	-
Gemini (Gemini Team et al., 2023)	-	-	-	-
Claude-3.5-Sonnet (Anthropic, 2024)		-	-	-
Deepseek-VL-Tiny (Lu et al., 2024)		SigLIP	DeepSeek LLM	MLP
Kimi-VL (Team et al., 2025)		MoonViT	Moonlight model	MLP
Llama-3.2-Vision-11B (Meta, 2024)		XAttn LLM	Llama 3.1	XAttn LLM
MiniCPM-0 2.6 (Yao et al., 2024)	× 7D	SigLIP	Qwen2.5-7B	MLP
Idefics3-8B-Llama3 (Laurençon et al., 2024)	> /B	SigLIP	Llama-3.1-8B-Instruct	XAttn LLM
InternVL2-8B (Chen et al., 2024d)		InternViT	internLM2.5-7b-chat	MLP
InternVL2.5-8B-MPO (Wang et al., 2024b)		InternViT-V2.5	internLM2.5-7b-chat	MLP
InternVL2.5-8B-MPO-AWQ (Chen et al., 2024c)		InternViT-V2.5	internLM2.5-7b-chat	MLP
LLaVA-v1.6-vicuna-7b (Liu et al., 2024)		CLIP ViT-L	vicuna-7b-v1.5	MLP
LLaVA-v1.6-mistral-7b (Liu et al., 2024)		CLIP ViT-L	Mistral-7B-Instruct-v0.2	MLP
Qwen2.5-VL-3B-Instruct (Team, 2025)		ViT	Qwen2.5 LLM	MLP
Qwen2.5-VL-7B-Instruct (Team, 2025)	$\leq 7B$	ViT	Qwen2.5 LLM	MLP
mPLUG-Owl3-7B-240728 (Ye et al., 2024)		SigLIP	Qwen2 LLM	Linear
mPLUG-Owl3-2B-241014 (Ye et al., 2024)		SigLIP	Qwen2 LLM	Linear
LLaVA-OneVision (Li et al., 2024)		SigLIP	Qwen2 LLM	MLP

Table 7: API-based and open-source MLLMs selected for evaluation.

Model	Polysemy		Homonymy		Grammar		Semantics		Specialized		Cultural		Overall	
Widdei	Accp	Accq	Accp	Accq	Accp	Accq	Accp	Accq	Acc _p	Accq	Acc _p	Accq	Accp	Accq
Qwen2.5-VL-7B-Instruct InternVL2.5-8B-MPO-AWQ LLaVA-OneVision MiniCPM-0 2.6	28.57 41.07 24.72 20.22	56.34 68.44 55.31 55.31	28.71 44.55 29.91 40.17	56.88 68.35 54.00 66.40	17.65 30.88 17.59 34.25	49.26 63.24 50.93 64.49	19.05 33.33 23.81 33.33	38.10 64.29 50.00 60.71	25.00 37.50 22.73 31.82	55.17 59.77 47.37 61.05	17.07 24.39 14.63 14.63	50.59 60.00 51.76 45.88	25.06 38.04 22.78 30.52	53.91 65.82 52.26 60.75

Table 8: We input prompts in Chinese, and the options are without confusing elements.

options enabled.

C More Cases

Model	Polys Acc _p	semy Acc _q	Homo Acc _p	onymy Acc _q	Gran Acc _p	nmar Acc _q	Sema Acc _p	ntics Acc _q	Specia Acc _p	alized Acc _q	Cult Acc _p	ural Acc _q	Ove Acc _p	erall Acc _q
Qwen2.5-VL-7B-Instruct InternVL2.5-8B-MPO-AWQ LLaVA-OneVision MiniCPM o 2.6	12.36 29.21 7.87	42.46 60.89 44.13	17.95 29.91 23.93 27.35	50.40 58.40 50.80	20.37 26.85 12.96	51.40 56.54 50.93	21.43 26.19 11.90 28.57	44.05 58.33 44.05	25.00 31.82 13.64	51.58 57.89 43.16	12.20 21.95 7.32	41.18 55.29 45.88	17.77 28.25 14.35	47.74 58.10 47.63

Table 9: We input prompts in Chinese, and the options are with confusing elements.

Model	Poly Acc _p	semy Acc _q	Homo Acc _p	onymy Acc _q	Gran Acc _p	nmar Acc _q	Sema Acc _p	antics Acc _q	Speci Acc _p	alized Acc _q	Cult Acc _p	ural Acc _q	Ove	erall Acc _q
Qwen2.5-VL-7B-Instruct	35.71	63.13	30.69	61.01	25.00	55.88	28.57	54.76	37.50	62.07	9.76	51.76	30.30	59.98
InternVL2.5-8B-MPO-AWQ	39.88	68.14	36.63	66.06	26.47	62.50	28.57	61.90	35.00	60.92	26.83	62.35	34.85	65.27
LLaVA-OneVision	20.22	20.28	25.64	56.00	17.60	49.53	26.19	53.57	15.91	46.32	19.51	50.58	20.73	51.60
MiniCPM-o 2.6	33.71	52.60	41.03	66.80	37.96	64.95	28.57	55.95	27.27	62.11	14.63	48.24	33.49	62.29

Table 10: We input prompts in Malay, and the options are without confusing elements.

Model	Poly: Acc _p	semy Acc _q	Homo Acc _p	onymy Acc _q	Gran Acc _p	nmar Acc _q	Sema Acc _p	ntics Acc _q	Specia Acc _p	alized Acc _q	Cult Acc _p	ural Acc _q	Ove	erall Acc _q
Qwen2.5-VL-7B-Instruct	14.61	45.25	18.80	50.00	24.07	50.00	21.43	46.43	25.00	51.58	9.76	36.47	19.13	47.63
InternVL2.5-8B-MPO-AWQ	23.60	53.07	21.37	53.20	24.07	55.14	23.81	57.14	27.27	55.79	12.20	44.88	22.32	53.58
LLaVA-OneVision	15.73	46.93	27.35	56.00	15.74	51.41	21.43	52.38	1.37	40.00	12.20	51.76	18.45	50.72
MiniCPM-o 2.6	22.47	48.60	30.77	59.60	37.04	61.68	19.05	44.05	18.18	53.68	9.76	36.47	26.20	53.69

Table 11: We input prompts in Malay, and the options are with confusing elements.

Q: 'Angkat bendera putih' merujuk kepada apa? (What does "raised white flag" mean?) C: Pasukan itu angkat bendera putih.

Scenario: 1



Answer: A.Mengaku kalah atau menyerah diri. (Admit defeat or surrender.)

Scenario: 1

Explanation: "Angkat bendera putih" is a widely recognized idiom, particularly in a military or competitive context. It is an international symbol signifying surrender, admission of defeat, or giving up. This clearly implies that the team has admitted defeat or surrendered.

X InternVL2_5-8B-MPO-AWQ: B.

✓ InternVL2_5-8B-MPO-AWQ (with CoT): A.

The image shows a group of people in camouflage outfits, likely engaged in a paintball game. They are positioned behind sandbags and a fence, with one person holding a white flag. In the context of paintball or similar games, raising a white flag typically indicates surrender or the end of the game. This is a common signal to show that the team is willing to concede or stop playing.\n\nAnswer: A

✓ Qwen2.5-VL-7B-Instruct: A.

✓ Qwen2.5-VL-7B-Instruct (with CoT): A.

Dalam gambar ini, terlihat dua pasukan yang sedang bermain paintball di balik sandbag dan pagar besi. Salah satu pasukan mengangkat bendera putih. Dalam konteks permainan paintball, mengangkat bendera putih biasanya merupakan tindakan untuk menunjukkan bahwa mereka telah menyerah atau mengakui kekalahan.\n\n**Answer: A** Scenario: 2



Answer: B.Mengibarkan bendera berwarna putih secara fizikal (Physically raising a white flag.)

Scenario: 2

Explanation: "Angkat bendera putih" can be interpreted as the physical act of raising a flag that happens to be white. It simply describes the action of lifting a whitecolored piece of fabric on a pole, without any inherent implication of defeat or surrender.

X InternVL2_5-8B-MPO-AWQ: A. X InternVL2_5-8B-MPO-AWQ (with CoT): A.

The image shows a person holding up a white flag in an open field. The act of raising a white flag typically signifies surrender or a request for peace. This is a common symbol in various contexts, including sports, military, and other situations where a peaceful resolution is sought. The other option, B, refers to physically holding up a white flag, which is not the primary meaning of the phrase "angkat bendera putih."\n\nAnswer: A

★Qwen2.5-VL-7B-Instruct: A. ✓ Qwen2.5-VL-7B-Instruct (with CoT): B.

'Angkat bendera putih' dalam konteks ini merujuk kepada tindakan mengibarkan bendera berwarna putih secara fizikal. Ini adalah tindakan yang biasanya dilakukan untuk menunjukkan kekuatan atau semangat dalam suatu perayaan atau acara tertentu.n\n<mark>Answer: B</mark>

Figure 6: A case of the category of polysemy in Malay.

Q:他想表达什么?(What is he trying to express?) C:"我家热得快炸了"。

Scenario: 1



Answer: A.我家非常热。 (My house is so hot.)

Scenario: 1

Explanation: "我家热得快炸了" can be punctuated in the way of "我家热得 | 快炸了", which is an exaggerated technique to describe "My house is so hot!" in Chinese.

✓ InternVL2.5-8B-MPO-AWQ: A.

✓ InternVL2.5-8B-MPO-AWQ (with CoT): A. To determine what the person is expressing in the phrase \"我 家热得快炸了\", we need to analyze the context and the literal meaning of the words. The phrase is in Chinese, and it translates to \"My house is so hot it's about to explode.\" The literal translation suggests extreme heat, which is a common way to express that a place is very hot. The phrase does not specifically mention a heating system or a malfunction, so it is more likely to be a general statement about the temperature rather than a technical issue.\n\nAnswer: A

✓ Qwen2.5-VL-7B-Instruct: A.

✓ Qwen2.5-VL-7B-Instruct (with CoT): A. \"我家热得快炸了\"这句话的意思是"我家热得快要爆炸 了",这是一种夸张的说法,用来形容非常热。选项A"我 家非常热"与这句话的意思相符;选项B"我家的加热器发 生故障了"则不符合这句话的含义。\n\n<mark>Answer: A</mark>

✓ MiniCPM-o 2.6: A.

To answer the question, let's break down the elements in the image and the context provided:\n\n1. The image shows a person standing in a forested area, with a background of lush greenery and a stream.\n2. The person is wearing a white t-shirt that appears to be wet, suggesting they might have been in water or sweating,\n3. The person is holding their head with one hand, which could indicate discomfort or distress.\n\nNow, let's analyze the options:\n\nA. \"我家非常热。\" - This translates to \"My house is very hot.\"\nB. \"我家的加热器发生故障了。\" - This translates to \"My heater is broken.\"\n/nGiven the context of the image:\n- The person is in a natural, outdoor setting, which is typically cooler than an indoor environment.\n- The person's wet t-shirt suggests they might be sweating or have been in water, which could be uncomfortable in a hot environment.\n\nConsidering these points, the most likely reason for the person's discomfort is the heat, rather than a broken heater,\n\nTherefore, the correct answer is:\n\nAnswer: A

Scenario: 2

Answer: B.我家的加热器发生故障了。 (There is a malfunction with the heater.)

Scenario: 2

Explanation: "我家热得快炸了" can be punctuated in the way of "我家热得快 | 炸了". Here, "热得快" is a common electric heater. Therefore, this sentence can be interpreted as "My '热得快' exploded!"

X InternVL2.5-8B-MPO-AWQ: A.

★ InternVL2.5-8B-MPO-AWQ (with CoT): A. The image shows a bucket of water with steam rising from it, indicating that the water is hot. The phrase \"我家热得快炸 了\" translates to \"My home is so hot it's about to explode.\" This suggests that the person is experiencing extreme heat in their home. The presence of steam and the context of the phrase imply that the heat is intense, but it does not necessarily mean that the heating device itself is malfunctioning.\n\nAnswer: A

✓ Qwen2.5-VL-7B-Instruct: B.

★ Qwen2.5-VL-7B-Instruct (with CoT): A.
\"我家热得快炸了\"这句话的意思是"我家热得快要爆炸了",这是一种夸张的说法,用来形容非常热。根据图片内容,可以看到一个装有水的容器中插着一根黄色的加热棒,水面上有蒸汽冒出,说明正在加热。因此,这句话可能是在描述加热器正在工作,导致家里变得非常热。 \n\nAnswer: A

✓ MiniCPM-o 2.6: B.

To answer the question, let's break down the elements in the image and the context provided:\n\n1. The image shows a container with water and a device that appears to be a heating element.\n2. The device is

submerged in water, and there is steam rising from the water.\n3. The text in the image is in Chinese: \"我家热得快炸了\" which translates to \"My home's instant heater is about to explode.\"\n\nNow, let's analyze the options:\n\nA. \"我家非常热。\" - This translates to \"My home is very hot.\" This option is not directly related to the image, as the image shows a heating device, not a hot environment.\n\nB. \"我家的加热器发生故障 了。\" - This translates to \"My home's heater is malfunctioning.\" This option is more relevant to the image, as it suggests a problem with the heating device.\n\nGiven the context of the image and the text, the most appropriate answer is:\n\n**Answer: B**

Figure 7: A case of the category of grammar in Chinese.

Q: Apa yang dibeli? (What was bought?)

C: Saya beli nasi lemak sambal sotong. (I bought nasi lemak with squid sambal.)



Answer: A.Hidangan nasi lemak yang berlauk sambal sotong. (Nasi lemak dish with squid sambal..)

Scenario: 1

Explanation: "nasi lemak sambal sotong" is a very common way to describe a single, complete dish in Malay cuisine. In this context, "sambal sotong" functions as the specific side dish or accompaniment that comes with the nasi lemak.

X InternVL2_5-8B-MPO-AWQ: B. X InternVL2 5-8B-MPO-AWQ (with CoT): B.

The image shows a plate of food with various components, including rice, boiled egg, cucumber slices, and two different types of side dishes. Based on the image, it is clear that the food items shown are nasi lemak and sambal sotong.\n6. Therefore, the correct option is B, which states that the person bought nasi lemak and sambal sotong.\n\Answer: B

✓ Qwen2.5-VL-7B-Instruct: A.

✔ Qwen2.5-VL-7B-Instruct (with CoT): A.

The image shows a plate of nasi lemak, which is a traditional Malaysian dish consisting of rice cooked in coconut milk and pandan leaves, served with various side dishes. The image clearly shows a plate of nasi lemak with sambal sotong as one of the side dishes. Therefore, the correct interpretation is that the person bought a plate of nasi lemak with sambal sotong.\n\n**Answer: A**



Answer: B.Nasi lemak dan juga sambal sotong. (Nasi lemak and squid sambal separately.)

Scenario: 2

Explanation: "nasi lemak sambal sotong" can be interpreted as two separate items: "nasi lemak" (a dish of coconut milk rice) and "sambal sotong" (a side dish of squid cooked in chili paste) being bought individually.

✓ InternVL2_5-8B-MPO-AWQ: B.

✓ InternVL2_5-8B-MPO-AWQ (with CoT): B.

The image shows a bowl of white rice, a bowl of dark sauce, and a bowl of what appears to be toasted bread. There is no indication of sambal sotong being a separate dish or item. The presence of the sauce and the bread suggests that the items are part of a meal, likely nasi lemak with sambal sotong as a condiment or side.\n\n**Answer: B**

XQwen2.5-VL-7B-Instruct: A. XQwen2.5-VL-7B-Instruct (with CoT): A.

The image shows a bowl of rice, a bowl of what appears to be sambal sotong (a type of chili paste), and some bread slices with the same chili paste spread on them. The presence of the rice and the chili paste strongly indicates that the dish being referred to is nasi lemak with sambal sotong as one of its components.n\nAnswer: A

Figure 8: A case of the category of grammar in Malay.

Q: Siapakah yang berada di bank? (Who was at the bank?)C: Dia nampak pencuri di bank. (He saw a thief at the bank.)



Answer: A.Pencuri itu berada di dalam atau di kawasan bank. (The thief was in or around the bank.)

♦ Scenario: 1

Explanation: "Dia nampak pencuri di bank" directly states that the thief was located at the bank. The most direct answer, based on the explicit mention, is the thief. This interpretation focuses on the stated location of the thief.

✓ InternVL2_5-8B-MPO-AWQ: A.

✓ InternVL2_5-8B-MPO-AWQ (with CoT): A.

To determine who is at the bank, we need to analyze the image carefully. The image shows a person dressed in black, wearing a mask and gloves, and appears to be attempting to steal from an ATM. This person is positioned in front of the ATM, indicating that they are physically present at the bank. The presence of the ATM and the person's actions suggest that they are inside or in the vicinity of the bank.\n\nAnswer: A

✔ Qwen2.5-VL-7B-Instruct: A.

✓ Qwen2.5-VL-7B-Instruct (with CoT): A.

The image shows a thief at an ATM machine, which is typically located inside or near a bank. The thief is depicted in a way that suggests he is inside the bank premises, as he is using the ATM. There is no indication that the person observing the thief is also inside the bank. Therefore, the correct interpretation is that the thief is inside the bank.\n\nAnswer: A

Scenario: 2



Answer: B.Dia berada di bank semasa melihat pencuri. (He was at the bank while seeing the thief.)

Scenario: 2

Explanation: "Dia nampak pencuri di bank" implies the observer's location. For "He saw a thief at the bank" (to occur, the pronoun "He" must also have been present "at the bank" to perform the act of seeing the thief there.

InternVL2_5-8B-MPO-AWQ: B. InternVL2_5-8B-MPO-AWQ (with CoT): B.

To determine who was at the bank, we need to analyze the image carefully. The image shows a person dressed in dark clothing, running down a street. However, there is no direct indication that the person in the hoodie was at the bank or that they were involved in a bank robbery. The image does not provide any clear evidence of the person being inside the bank or at the bank during the time of the alleged robbery."\n\nAnswer: A

X Qwen2.5-VL-7B-Instruct: B.

✔ Qwen2.5-VL-7B-Instruct (with CoT): A.

The image shows a person wearing a black hoodie and gloves, running away from the camera. The person in the foreground is not inside the bank but rather running away from it. Therefore, the correct interpretation is that the person is not inside the bank but rather outside, possibly after committing a crime.n\nAnswer: A

Figure 9: A case of the category of grammar in Malay.

Q: What is the subject in the sentence going to eat? C: The chicken is ready to eat.



Answer: A.Chicken.

♦ Scenario: 1

Explanation: The chicken is cooked and prepared, so it is ready for someone to eat.

✓ InternVL2.5-8B-MPO-AWQ: A.

✓ InternVL2.5-8B-MPO-AWQ (with CoT): A. The image shows a chicken inside a pot, which suggests that the chicken is being cooked. The sentence \"The chicken is ready to eat/" implies that the chicken has been prepared for consumption. Therefore, the subject in the sentence is likely to be the person or entity that is going to eat the <mark>chicken.</mark>\n\n<mark>Answer: A</mark>

✓ Qwen2.5-VL-7B-Instruct: A.

✓ Qwen2.5-VL-7B-Instruct (with CoT): A.

The image shows a cooked chicken on a plate with steam rising from it, indicating that it is hot and ready to be eaten. The question asks what the subject in the sentence is going to eat. Since the subject is implied to be a person who would eat the chicken, the correct answer should refer to the chicken itself as the food item being consumed.\n\nAnswer: A

✓ MiniCPM-o 2.6: A.

We found that MiniCPM-0 2.6's instruction-following ability isn't very good, so it didn't provide an explanation here. Scenario: 2



Answer: B.Chicken feed.

♦ Scenario: 2

Explanation: The chicken itself is hungry and ready to eat something.

X InternVL2.5-8B-MPO-AWQ: A. X InternVL2.5-8B-MPO-AWQ (with CoT): A.

The image shows a chicken sitting at a table with a plate, fork, and knife in front of it. The chicken appears to be ready to eat, suggesting that it is the subject of the sentence. The options provided are \"A. Chicken.\" and \"B. Chicken feed.\" Since the chicken is the one that is ready to eat, the subject in the sentence is the chicken itself.\n\nAnswer: A

X Qwen2.5-VL-7B-Instruct: A. X Qwen2.5-VL-7B-Instruct (with CoT): A.

The image shows a cartoon chicken sitting at a table with a plate and utensils, suggesting that the chicken is about to eat. The question asks what the subject (the chicken) is going to eat. Since there is no indication of the chicken eating anything other than food, the logical conclusion is that the chicken is going to eat something.\n\nAnswer: A

X MiniCPM-o 2.6: A.

We found that MiniCPM-0 2.6's instruction-following ability isn't very good, so it didn't provide an explanation here.

Figure 10: A case of the category of semantics in English.



Scenario: 1

Explanation: In the sentence "山上有许多杜鹃", "杜鹃" can be interpreted as an animal "cuckoo" or a plant "azalea". The sentence "惆怅多山人复稀 杜鹃啼处泪沾衣" in the image comes from an ancient Chinese poem 忆故圆, "杜鹃" here means cuckoo.

✓ InternVL2.5-8B-MPO-AWQ: A.

✓ InternVL2.5-8B-MPO-AWQ (with CoT): A.

The image contains Chinese text, and the question is asking whether \"山上有许多杜鹃\" refers to plants or animals. The text in the image is a poem, and the phrase \"杜鹃\" in Chinese refers to a bird, specifically the cuckoo. Therefore, the correct interpretation is that it refers to an

animal.\n\n<mark>Answer: A</mark>

X Qwen2.5-VL-7B-Instruct: B.

✔ Qwen2.5-VL-7B-Instruct (with CoT): A.

首先,我们需要理解图片中的文字内容。图片中的文字是 "惆怅多山人复稀 杜鹃啼处泪沾衣",这是一首古诗的诗 句。<mark>诗句中提到"杜鹃",在中国文化中,"杜鹃"通常指的</mark> 是杜鹃鸟,这是一种鸟类。因此,根据诗句的内容,我们 可以推断出"山上有许多杜鹃"是指杜鹃鸟而不是植物。 \n\n<mark>Answer: A</mark>

✓ MiniCPM-o 2.6: A.

To determine whether \"山上有许多杜鹃\" refers to plants or animals, we need to understand the meaning of the phrase in Chinese. \n\n\"杜鹃\" (dù quān) refers to a type of bird, specifically the Asian koel, which is known for its loud and distinctive call. The phrase \"山上有许多杜鹃\" translates to \"There are many Asian koels on the mountain.\"\n\nGiven this translation, the correct answer is that \"杜鹃\" refers to animals, specifically birds.\n\nTherefore, the correct option is:\n\n<mark>Answer: A</mark>

Scenario: 2

Explanation: In the sentence "山上有许多杜鹃", "杜鹃" can be interpreted as an animal "cuckoo" or a plant "azalea". The image shows a patch of flowers. "杜鹃" here means azalea.

✓ InternVL2.5-8B-MPO-AWQ: B.

InternVL2.5-8B-MPO-AWQ (with CoT): B. The image shows a hillside covered with vibrant pink

flowers. These flowers are typically associated with the plant species Rhododendron, commonly known as azaleas or rhododendrons. The term \"杜鹃\" in Chinese refers to these flowering plants, not animals. Therefore, the correct interpretation of \"山上有许多杜鹃\" in the context of the image is that there are many plants on the mountain.\n\nAnswer: B

✔ Qwen2.5-VL-7B-Instruct: B.

✓ Qwen2.5-VL-7B-Instruct (with CoT): B. 根据图片内容,可以看到山上有很多红色的花朵,这些 花朵是杜鹃花。杜鹃花是一种植物,因此这句话指的是 植物。\n\n<mark>Answer: B</mark>

✔ MiniCPM-o 2.6: B.

To determine whether \"山上有许多杜鹃\" refers to plants or animals, we need to analyze the context and the image provided.\n\n1. **Context Analysis**:\n - The phrase \"山 上有许多杜鹃\" translates to \"There are many azaleas on the mountain.\"\n - Azaleas are a type of flowering plant, not an animal.\n\n2. **Image Analysis**:\n - The image shows a mountainous area with a lot of greenery and vibrant pink flowers.\n - The pink flowers are clearly visible and appear to be azaleas, which are plants.\n\n3. **Options Analysis**:\n - Option A: 动物 (Animal) - This does not fit the context as azaleas are not animals.\n - Option B: 植物 (Plant) - This fits the context as azaleas are plants.\n\nBased on the context and the image, the correct answer is that \"山 上有许多杜鹃\" refers to plants.\n\nAnswer: B

Figure 11: A case of the category of specialized in Chinese.

Q: The cultural background of the customer is shown in the picture, what will the customer's attitude be?

C: "Please tip me".



Answer: A.Positive.

Scenario: 1

Explanation: The picture shows New York City in the United States. There is a tipping culture in Western societies. The tipping culture is something they are quite accustomed to.

✓ InternVL2.5-8B-MPO-AWQ: A.

✓ InternVL2.5-8B-MPO-AWQ (with CoT): A. To answer the question, consider the following steps:\n\n1. The image shows a cityscape with a prominent skyline, which is likely New York City.\n2. New York City is known for its diverse cultural background and multicultural environment. 3. The phrase \"Please tip me\" is a common request in service industries, such as restaurants or bars.\n4. In New York City, tipping is generally expected in these settings as a part of the service culture.\n5. Given the cultural context, the customer's attitude is likely to be positive, as they are familiar with and expect tipping in such situations.\n\n**Answer: A**

✓ Qwen2.5-VL-7B-Instruct: A. XQwen2.5-VL-7B-Instruct (with CoT): B.

The image depicts a city skyline at dusk with illuminated buildings, suggesting a modern urban environment. There is no direct indication of the cultural background of the customer or their attitude towards tipping. The question seems to be asking about the cultural implications of tipping, but the image does not provide enough context to determine the customer's attitude.\n\n**Answer: B**

? MiniCPM-o 2.6:

I'm sorry, but I can't provide assistance with that request.



Answer: B.Confused.

Scenario: 2

Explanation: The picture shows Chongqing City, China. There is no tipping culture in Eastern societies. Chinese people may feel confused about the act of asking for tips.

X InternVL2.5-8B-MPO-AWQ: A.

✓ InternVL2.5-8B-MPO-AWQ (with CoT): B. To answer the question, consider the following steps:\n\n1. **Cultural Context**: The image shows a traditional Chinese building illuminated at night, with a bridge and a river in the foreground. This setting is likely in a city with a rich cultural heritage, possibly in China.\n\n2. **Understanding the Phrase**: The phrase \"Please tip me\" is a request for a gratuity, which is a common practice in many Western cultures. Tipping is not a traditional practice in China, and it can sometimes be misunderstood or seen as unnecessary,\n\n3. **Cultural Sensitivity**: Given the cultural background of the customer, who is likely from a region where tipping is not customary, the customer might be confused or unfamiliar with the concept of tipping.\n\n4. **Possible Attitudes**: Based on the cultural context, the customer might be confused about why they are being asked to tip, as it is not a common practice in their culture.\n\nAnswer: B

X Qwen2.5-VL-7B-Instruct: A. X Qwen2.5-VL-7B-Instruct (with CoT): A.

The image depicts a vibrant cityscape at night with illuminated buildings and a bridge. The architecture suggests a blend of traditional and modern elements, possibly indicating a rich cultural heritage. The presence of people gathered along the riverbank adds to the lively atmosphere. Given the context of the image, which showcases a bustling urban environment with a mix of historical and contemporary influences, it can be inferred that the customer might have a positive attitude towards the cultural background depicted.\n\n**Answer: A**

? MiniCPM-o 2.6:

I'm sorry, but I can't provide assistance with that request.

Figure 12: A case of the category of cultural in English.





boat.

	Q: Who is my wife, the wo	man on the left or the right?
Scenario: 1	Scenario: 2	Scenario: 3
C: I'm <u>70</u> years old now. My ife and I met a couple on the	C: I'm <u>35</u> years old now. My wife and I met a couple on the boat	C: I'm <u>70</u> years old now. My wife and I met a couple on the boat
Answer: A.On the right.	Answer: B.On the left.	Answer: B.On the left.

◆ Scenario: 1

w

b

Explanation: 70 indicates that the speaker is elderly, and his wife is likely of similar age. This rules out the younger couple. The phrase "on the boat" is disambiguated to refer to the couple they met, not themselves. Based on the image, the elderly woman on the right is inferred to be his wife.

✓ InternVL2.5-8B-MPO-AWQ: A. XInternVL2.5-8B-MPO-AWQ (with CoT): B.

To determine who the woman on the left or the right is, 1 will analyze the image and look for clues that can help identify the woman in question. The question states that the person asking the question is 70 years old and is asking about their wife. The woman on the left is warning a wedding dress, which typically indicates that she is a bride. **The** woman on the right is an older couple, which does not match the description of the person asking the question being 70 years old. Therefore, the woman on the left is more likely to be the wife of the person asking the question. **Answer. B**

XQwen2.5-VL-7B-Instruct: B. ✓ Qwen2.5-VL-7B-Instruct (with CoT): A.

The image shows two older individuals standing in the foreground, likely the person asking the question and their wife. In the background, there is a couple on a boat. Given that the person asking the question is 70 years old and they met someone on the boat, it is reasonable to infer that the couple on the boat consists of the person's wife and another individual. Since the person asking the question is older and standing next to the woman on the right, it is logical to conclude that the woman on the right, is their wife Answer: A

♦ Scenario: 2

Explanation: 35 indicates that the speaker is relatively young, so his wife is likely also young. In this context, "on the boat" is interpreted to include the speaker and his wife. Given that the younger couple is in the boat, the woman on the left is inferred to be his wife.

✓ InternVL2.5-8B-MPO-AWQ: B. ✓ InternVL2.5-8B-MPO-AWQ (with CoT): B.

To determine who the woman on the left or the right is the speaker's wife, we need to consider the context provided. The speaker is 35 years old and is asking about his wife. Typically, in a wedding seenario, the bride is the woman who is married to the groom. In the image, the woman on the left is wearing a wedding dress and veil, which are traditional attire for a bride. The woman on the tright is older and dressed in a more casual manner, which is less likely to be the bride. Therefore, based on the traditional roles and appearances, the woman on the likely to be the speaker's wife. Amswer: **B**

Qwen2.5-VL-7B-Instruct: B. Qwen2.5-VL-7B-Instruct (with CoT): B.

The image shows two people in the foreground who appear to be older adults, likely grandparents, and two younger individuals in the background who seem to be a couple. The person on the left in the background is wearing a wedding dress, indicating that they are likely the bride. The person on the right in the background is dressed in formal attire, suggesting they might be the groom. Since you mentioned meeting a couple on the boat, the woman on the left is more likely to be your wife as the is wearing a wedding dress. Answer: B

♦ Scenario: 3

Explanation: 70 implies the speaker is elderly, so his wife is expected to be the elderly woman. The phrase "on the boat" refers to the couple they met. Given the cruise setting, and with visual support, the woman on the left (elderly) is identified as the speaker's wife.

InternVL2.5-8B-MPO-AWQ: B. InternVL2.5-8B-MPO-AWQ (with CoT): A.

To determine who the woman on the right is, I will analyze the image step by step. 1. The question states that the person asking the

question is 70 years old and is asking about their wife. 2. The image shows two couples on a boat.

3. The woman on the right is wearing a yellow dress and has long brown hair.

 The woman on the left is wearing a pink top and a straw hat.
 Typically, a person would refer to their spouse as

'my wife" when they are together.

6. Since the person asking the question is a man, and he is asking about his wife, it is likely that the woman on the right is his wife. Answer: A

✓ Qwen2.5-VL-7B-Instruct: B. ✓ Qwen2.5-VL-7B-Instruct (with CoT): B.

The image shows four people standing on a boat deck. There is an older man and woman on the left side of the image, and a younger man and woman on the right side. The older man and woman appear to be shaking hands with the younger couple. Given that the question states the person asking is 70 years old and met a couple on the boat, fit is reasonable to infer that the older man and woman are the person's wife and husband. Therefore, the woman on the left is the person's wife. <u>Answer: B</u>

Answer: A.On the right.

Scenario: 4

C: I'm 35 years old now. My

wife and I met a couple on the

Explanation: 35 suggests the speaker is young, making it likely that his wife is also from the younger couple. The phrase "on the boat" again includes the speaker and his wife. Based on the cruise ship setting, the younger woman on the right is inferred to be the wife.

InternVL2.5-8B-MPO-AWQ: A. InternVL2.5-8B-MPO-AWQ (with CoT): A.

To determine who the woman on the right is, we need to consider the context provided in the question. The question states that the person is 35 years old now and is asking about their wife. Since the person is currently 35, it implies that they are married to the woman who is younger than them. In the image, the woman on the right appears to be younger than the woman on the left. Therefore, the woman on the right is likely the younger woman, who would be the wife of the 35-year-old person. Answer: **A**

Qwen2.5-VL-7B-Instruct: A. Qwen2.5-VL-7B-Instruct (with CoT): A.

The image shows four people standing on a boat deck. There is an older couple on the left and a younger couple on the right. <u>Since your mentioner</u> that you are 35 years old and met a couple on the boat, the younger couple on the right would be more likely to be your age. Therefore, the woman on the right is probably your wife. <u>Answer: A</u>

Figure 13: A case of the category of dual-ambiguity in English.