Drawing the Line: Enhancing Trustworthiness of MLLMs Through the Power of Refusal

Anonymous ACL submission

Abstract

Multimodal large language models (MLLMs) excel at multimodal perception and understanding, yet their tendency to generate hallucinated or inaccurate responses undermines their trustworthiness. Existing methods have largely overlooked the importance of refusal responses as a means of enhancing MLLMs reliability. To bridge this gap, we present the Information Boundary-aware Learning Framework (InBoL), a novel approach that empowers MLLMs to refuse to answer user queries when encountering insufficient information. To the best of our knowledge, InBoL is the first framework that systematically defines the conditions under which refusal is appropriate for MLLMs using the concept of information boundaries proposed in our paper. This framework introduces a comprehensive data generation pipeline and tailored training strategies to improve the model's ability to deliver appropriate refusal responses. To evaluate the trustworthiness of MLLMs, we further propose a user-centric alignment goal along with corresponding metrics. Experimental results demonstrate a significant improvement in refusal accuracy without noticeably compromising the model's helpfulness, establishing InBoL as a pivotal advancement in building more trustworthy MLLMs.

1 Introduction

004

011

012

014

018

023

034

042

Recent advancements in multimodal large language models (MLLMs) have marked a significant breakthrough in AI research, especially in visionlanguage tasks (McKinzie et al., 2024; Bai et al., 2023; Tong et al., 2024; Fu et al., 2024; Li et al., 2024; Zhang et al., 2024b). By integrating visual information with large language models (LLMs), these models have exhibited profound capabilities in multimodal understanding and reasoning. Despite the impressive progress, MLLMs still face notable challenges. One prominent issue is their tendency to generate factually incorrect or hallucinated content, where models confidently describe non-existent visual elements or provide responses that include incorrect knowledge (Bai et al., 2024; Zhong et al., 2024). Such hallucinations not only reduce the accuracy of the models but also undermine their truthfulness in practical applications, hindering them from being trustworthy AI systems. 043

045

047

051

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

083

To improve the trustworthiness of MLLMs, previous works primarily focus on improving multimodal alignment algorithms to enhance the models' perceptual and reasoning capabilities (Yu et al., 2024a,b; Amirloo et al., 2024). However, all models have intrinsic limitations in their knowledge and perceptual capabilities, making them prone to produce inaccurate or misleading responses when confronted with tasks beyond their capabilities. Therefore, another effective approach to improving trustworthiness is to train these models to recognize their boundaries and refuse to answer questions when appropriate, thus preventing the generation of misinformation.

Despite the critical role of refusal responses, few studies have focused on effectively training MLLMs for this capability. Existing approach (Cha et al., 2024) primarily targets ambiguous or unanswerable queries, such as those involving nonexistent visual elements, but fall short of addressing the broader challenges related to intrinsic limitations and self-awareness in MLLMs (Wang et al., 2024b). This gap underscores the need for strategies that enable MLLMs to recognize their limitations, ensuring they either provide accurate responses or appropriately refuse to answer when necessary.

While related research in MLLMs is limited, efforts to improve reliability by training models to refuse answering unknown questions have been widely studied in LLMs (Amayuelas et al., 2024; Yin et al., 2023; Yang et al., 2023; Cheng et al., 2024; Chen et al., 2024; Liang et al., 2024). These studies typically generate instruction and preference data that include refusal responses to guide models in avoiding responses beyond their knowledge boundaries. However, extending these approaches to MLLMs presents unique challenges. Unlike LLMs, the trustworthiness of MLLMs relies not only on internal knowledge but also on external multimodal inputs, which are unaccounted for in existing methods. This raises the critical question of when MLLMs should refuse to answer. Furthermore, evaluating trustworthiness in these works involves categorizing test questions as known and unknown based on knowledge boundaries—a process that requires high computational cost as it demands multiple sampling.

084

086

090

100

101

102

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

129

130

131

132

133

134

To address these limitations, we introduce the Information Boundary-aware Learning Framework (InBoL), the first to establish the concept of information boundaries and systematically define the conditions under which MLLMs should appropriately refuse to respond. This marks a significant advancement in trustworthiness training for MLLMs. Building on these boundaries, we develop a data construction pipeline that generates "I Don't Know" (IDK) instruction and preference data from available VQA datasets. Using this data, we implement two training methods: IDK Instruction Tuning (IDK-IT) and Confidence-aware Direct Preference Optimization (CA-DPO), which enables MLLMs to recognize their information boundaries and refuse to answer when necessary. To evaluate model trustworthiness, we introducing simple and model-agnostic metrics for assessing trustworthiness from a user-centric perspective. Our experimental results demonstrate that InBoL significantly improves the trustworthiness of baseline models by enhancing their ability to appropriately refuse responses while maintaining helpfulness. This work introduces a new paradigm for developing trustworthy MLLMs and sets the foundation for future advancements in this critical area.

Overall, our contributions are as follows:

• InBoL Framework: We propose the InBoL framework, which introduces the novel concept of information boundaries and integrates a comprehensive data construction pipeline along with tailored training methods. InBoL enhances the trustworthiness of MLLMs by empowering them to recognize these boundaries and refuse to answer when lacking sufficient information, setting a new benchmark for trustworthiness training.

• User-centric Trustworthiness Evaluation:

We introduce a user-centered alignment objective and define several simple and model-
agnostic metrics for trustworthiness, which135simplifies the evaluation process.138

• Experimental Validation: We conduct extensive experiments to validate the effectiveness of our approach, demonstrating significant improvements in MLLMs' ability to recognize information boundaries while preserving helpfulness. Our detailed analyses offer valuable insights into the broader impact of this method, paving the way for future developments in trustworthy MLLMs.

2 Problem Formulation

2.1 MLLM Alignment for Trustworthiness

Previous studies have explored alignment for trust-150 worthiness in LLMs, primarily focusing on eval-151 uating trustworthiness based on the model's in-152 trinsic boundary. Given a user query q and the 153 model-generated response r, the trustworthiness 154 of the response is evaluated by a value function 155 $v(q,r) \in \{0,1\}$. The goal of alignment is to max-156 imize $\sum_{q \in D_{\text{test}}} v(q, r)$. To determine the value of 157 $v(\cdot)$, these studies first classify the test set ques-158 tions into known D_k and unknown D_{uk} categories 159 based on the model's knowledge boundary. The 160 value function $v(\cdot)$ is then defined as: 161

$$v(q,r) = \begin{cases} 1 & \text{if } q \in D_k \text{ and } r \text{ is correct.} \\ 1 & \text{if } q \in D_{uk} \text{ and } r \text{ is a refusal response.} \\ 0 & \text{otherwise} \end{cases}$$
(1)

162

139

140

141

142

143

144

145

146

147

148

149

However, categorizing questions as 'known' or 163 'unknown' for each model is a complex task with 164 high computational costs. Additionally, D_k and 165 D_{uk} are model-specific, making it difficult to fairly 166 compare the trustworthiness of different models. 167 To address these challenges, we propose a new 168 model-agnostic alignment objective for trustworthiness. Inspired by Xu et al. (2024), our ap-170 proach evaluates trustworthiness based on user pref-171 erences: users value accurate, relevant, and infor-172 mative responses; prefer refusals over misinforma-173 tion; and find incorrect answers highly harmful. 174 A trustworthy MLLM should maximize accurate 175 responses while using refusals to prevent misinfor-176 mation. Hence, we redefine the value function as 177 follows: 178



Figure 1: Information Boundaries of MLLMs. (a) Questions are categorized into three types based on intrinsic and extrinsic information boundaries. For Type 1 questions, which fall within the intrinsic boundary, the model is expected to provide helpful responses. For Type 2 questions, which require knowledge beyond the model's capabilities (e.g., dog breeds), the model should refuse to answer. For Type 3 questions, where the provided image lacks sufficient information, the model should also respond with a refusal. (b) The intrinsic and extrinsic boundaries are illustrated, highlighting the model's varying confidence in answering queries across different regions.

$$v(q,r) = \begin{cases} 1 & \text{if } r \text{ is a correct response,} \\ 0 & \text{if } r \text{ is a refusal response,} \\ -1 & \text{if } r \text{ is a incorrect response.} \end{cases}$$
(2)

Consequently, the new objective for trustworthiness alignment is to maximize the sum of values over the test set:

$$\underset{\theta}{\text{maximize}} \quad \sum_{q \in D_{\text{test}}} v(q, r) \tag{3}$$

This objective encourages models to generate as many correct responses as possible while prioritizing refusal when accuracy cannot be guaranteed. Unlike previous approaches, the definition of $v(\cdot)$ is model-agnostic, allowing for consistent evaluation across different models.

2.2 Evaluation Metrics

179

182

184

185

188

189

190

191

193

198

200

To evaluate the model's trustworthiness, we intoduce two key metrics—Accuracy (Acc) and Refusal Rate (RefR)—defined as follows:

$$Acc = \frac{N_c}{N}, \quad \text{RefR} = \frac{N_r}{N} \tag{4}$$

where N_c is the number of correct responses, N_r is the number of refusal responses, and N is the total number of queries.

Combining these two metrics, we define the objective for trustworthiness alignment as trustworthiness score s_{trust} as follows:

$$s_{\text{trust}} = \sum_{q \in D_{\text{test}}} v(q, r) = 2 \cdot \text{Acc} + \text{RefR} - 1.$$
(5)

201

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

221

222

225

226

3 Information Boundary-Aware Learning framework

To enhance the trustworthiness of MLLMs, we propose the **In**formation **Bo**undary-Aware Learning Framework (InBoL). This framework includes a data construction pipeline designed to generate model-specific 'IDK' instruction and preference data by considering the intrinsic and extrinsic information boundaries. Furthermore, we incorporate IDK-IT and CA-DPO for model training. The goal of this framework is to improve the model's ability to provide appropriate refusal responses while maintaining helpfulness, thereby reducing misinformation and increasing reliability of MLLMs.

3.1 Information Boundary

The core of our framework is to train MLLMs to recognize when to refuse, thereby avoiding the generation of misinformation. While previous work on LLMs generally restricts refusals to questions outside the model's knowledge boundary, multimodal scenarios introduce additional complexity, as both visual information and knowledge must be considered.

To address this, we introduce extrinsic and intrinsic information boundaries for MLLMs, as il-



Figure 2: The Pipeline of Data Construction: Given a VQA dataset, we design a pipeline to collect different types of samples within and beyond the information boundaries. First, we estimate the confidence for each sample to determine the model's intrinsic information boundary. Next, we generate questions that lie beyond the extrinsic boundary, followed by quality filtering. Finally, all data is formatted into a standardized structure, including correct, incorrect, and refusal responses, each accompanied by their corresponding confidence scores.

 lustrated in Figure 1. In multimodal scenarios, a trustworthy MLLM should answer questions only when it has sufficient information and refuse when it does not, and these boundaries serve as guidelines for this decision-making process.

227

231

232

233

240

254

Extrinsic Information Boundary In multimodal scenarios, MLLMs depend on extrinsic visual inputs to respond to user queries. The extrinsic information boundary defines the distinction between what is explicitly present in the visual input and what is absent. When the necessary information to answer a question is not available—indicating that the query exceeds the extrinsic boundary—the model should provide a refusal response.

Intrinsic infomration boundary Beyond the ex-241 trinsic boundary, a model's intrinsic information 242 boundary is equally important, defined by its in-243 herent capabilities. This boundary encompasses what the model can infer from the image and the 246 multimodal knowledge embedded in its parameters. If the model cannot perceive the required informa-247 tion from the image or lacks specific knowledge, thereby exceeding the intrinsic boundary, it should 249 also provide a refusal response. 250

3.2 Data Construction

To train MLLMs to appropriately refuse questions, we need to collect the three types of VQA data outlined in Figure 1. For any given VQA dataset, we propose a data construction pipeline that classifies questions into the first two types based on confidence estimation, while generating unanswerable questions as the third type from the available data. Additionally, we reorganize the generated data into a standardized format, as shown in Figure 2, which is then used to create the 'IDK' instructions and preference data. 255

256

257

259

261

262

263

265

266

267

269

270

271

272

273

274

275

276

277

278

279

281

283

Estimating the Model's Confidence We assume that all questions in a given VQA dataset are answerable based on the provided visual information, placing them within the extrinsic information boundary. To further determine if these questions fall within the intrinsic information boundary, we estimate the model's confidence. Following prior works (Cheng et al., 2024; Xu et al., 2024; Yang et al., 2023), we sample ten responses from the original model and calculate the accuracy rate to estimate its confidence. In addition, we randomly sample a correct and incorrect answers from the model-generated responses and select a refusal response from the refusal template (Appendix B), then restructure the data into a standardized format as illustrated in Figure 2. Further details about this process are provided in Appendix C.

Generating Unanswerable Questions Next, we generate questions that lie beyond the extrinsic information boundary. First, we focus on questions that are irrelevant to the provided image. To cre-



Figure 3: Construction of 'IDK' instruction and preference data: The restructured data is categorized into 'Known,' 'Mixed,' and 'Unknown' based on confidence thresholds(δ_k and δ_{uk}). 'IDK' instruction generation includes correct responses for known questions, refusal responses for unknown questions, and the exclusion of mixed data. Preference data samples are constructed by pairing questions with correct, incorrect, and refusal responses, based on the confidence classification of each question.

ate these, we randomly select samples from the VQA dataset and reorganize them into mismatched image-question pairs. Additionally, we formulate more complex questions that, while related to the image, cannot be answered due to incorrect assumptions or insufficient information provided by the image. To achieve this, we design prompts that instruct GPT-40 to generate unanswerable queries and further implement a filtering mechanism to ensure the quality of them. For each generated unanswerable question, we assign a confidence score of 0. Additionally, we collect an incorrect response generated by the original model and select a refusal response from predefined templates to restructure the data into the standardized format shown in Figure 2. Detailed descriptions of the generation and filtering processes are provided in Appendix C.2 and we further discuss the generalization of this data generation method in Appendix F.3.

289

290

293

294

303Constructing 'IDK' InstructionTo construct304the 'IDK' instruction, we categorize the restruc-305tured data into three types based on the confidence306thresholds δ_k and δ_{uk} : 'Known,' 'Mixed,' and 'Un-307known,' as shown in Figure 3. For known questions,308we select the correct answer as the response. For309unknown questions, we utilize the refusal response.310Regarding the 'Mixed' data, we exclude it from the311instruction data, as the model exhibits relatively312high uncertainty for these questions.

313Constructing Preference DataA preference314data consists of a question, a chosen response, and315a rejected response. For known questions, we use316the correct answers as the chosen response and the

refusal as the rejected response. For unknown questions, we utilize the refusal answers as the chosen response and the incorrect answers as the rejected response. For mixed questions, we construct two samples, both of which use incorrect answers as the rejected response. In one sample, the chosen response is the correct answer, while in the other, the chosen response is the refusal. 317

318

319

320

321

322

323

324

325

327

329

331

332

333

334

335

337

338

340

341

342

343

344

345

347

348

3.3 Model Training for Information Boundary Awareness

To enhance the model's ability to recognize and refuse questions beyond its information boundary, we propose two training strategies: 'IDK' Instruction Tuning (IDK-IT) and Confidence-aware Direct Preference Optimization (CA-DPO). These strategies can teach the model to provide refusal responses when it lacks the necessary information.

IDK Instruction Tuning Instruction tuning is an effective method for aligning the model's responses with desired behaviors. In our framework, we train the model with 'IDK' instructions. This training approach improves models trustworthiness by reducing the generating misinformation.

Confidence-aware Direct Preference Optimization Direct Preference Optimization (DPO) is a technique that optimizes a model's policy using preference data (Xu et al., 2023; Rafailov et al., 2024; Hong et al., 2024; Yuan et al., 2024). While DPO can guide models to prefer correct answers and learn to refuse when needed, it does not leverage the model's intrinsic confidence to dynamically adjust its behavior. To address this, we pro-

Mathad	AOKVQA				GQA			MMMU			BeyondVisQA MMBench		Bench(e	(en-dev)	
Methou	Acc	RefR	S_{trust}	Acc	RefR	S_{trust}	Ā	Acc	RefR	S_{trust}	RefR	Acc	RefR	S_{trust}	
LLaVA1.5-7B	78.56	0.00	57.13	59.65	0.00	19.30	34	4.70	0.00	-30.60	25.50	62.80	0.00	25.60	
+Refusal Prompt	56.77	26.20	39.74	58.65	3.43	20.74	32	2.22	12.89	-22.67	27.50	59.36	0.69	19.42	
+SFT	74.32	3.49	52.14	59.39	2.77	21.55	34	4.20	1.67	-29.93	56.00	63.32	0.26	26.89	
+IDK-IT	55.50	36.24	47.24	50.46	23.88	24.81	1:	5.22	69.67	0.11	75.25	46.39	39.09	31.87	
+CA-DPO	72.23	17.64	62.10	60.41	12.95	33.77	19	9.67	56.67	-4.00	67.75	58.42	18.13	34.97	
LLaVA1.5-13B	78.95	0.00	57.90	61.81	0.00	23.63	30	6.22	0.00	-27.56	33.50	67.96	0.00	35.91	
+Refusal Prompt	63.32	18.95	45.59	61.36	1.96	24.69	2	7.78	19.56	-24.89	46.00	64.69	0.26	29.64	
+SFT	77.82	2.62	58.25	61.32	1.69	24.33	3	8.22	1.78	-21.78	68.75	67.01	0.00	34.02	
+IDK-IT	63.93	23.06	50.92	52.27	19.22	23.77	14	4.22	74.33	2.78	79.50	55.84	23.91	35.60	
+CA-DPO	73.89	15.63	63.41	59.70	13.82	33.22	2:	5.89	41.78	-6.44	72.50	62.63	14.69	39.95	

Table 1: Performance on out-of-domain dataset. We present results on LLaVA1.5-7B and LLaVA1.5-13B. Bold values indicate the highest trustworthiness score.

pose Confidence-aware DPO (CA-DPO), which integrates the model's confidence score into the optimization process. As shown in Figure 3, we define two preference pairs for 'Mixed' samples: p_1 (correct > incorrect) and p_2 (refusal > incorrect). For consistency, we define 'Known' samples with $p_1 = p_2 =$ (correct > refusal), and 'Unknown' samples with $p_1 = p_2 =$ (refusal > incorrect). Our approach uses the confidence score to dynamically balance the emphasis between these preference pairs. The CA-DPO loss function is defined as:

354

355

357

363

$$f(x,p) = \log \sigma(\beta \log \frac{\pi_*(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_*(y_l|x)}{\pi_{\text{ref}}(y_l|x)}$$
(6)

$$\mathcal{L}_{\text{cadpo}} = -\mathbb{E}_{(x,p_1,p_2)} \Big(f(x,p_1) \cdot conf_x + f(x,p_2) \cdot (1-conf_x) \Big)$$
(7)

where $conf_x$ is the model's confidence score. The confidence score adjusts the balance between the two preference pairs, particularly for 'Mixed' samples. In high-confidence scenarios, the loss function prioritizes correct responses, while in lowconfidence cases, it favors refusal. This adaptive mechanism enables the model to balance cautiousness and helpfulness more effectively.

4 Experimental Setup

Training Data As mentioned in Section 3.1, our work considers both the model's knowledge and visual information. Therefore, we use general VQA datasets and knowledge-intensive VQA datasets for data construction. Specifically, we utilize VQAV2 (Antol et al., 2015; Zhang et al., 2016; Goyal et al., 2017), Oven (Hu et al., 2023), and ScienceQA (Lu et al., 2022). We set the confidence thresholds as $\delta_k = 0.8$ and $\delta_{uk} = 0.2$. For 'IDK' instruction tuning, we collected 11k instructions. For CA-DPO, we gathered about 24k preference pairs. Further training details are provided in Appendix D.

383

384

385

387

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

Evaluation In our experiments, we utilize the LLaVA1.5 (Liu et al., 2023c,b) model, one of the most widely used open-source MLLMs. We evaluate models on both in-domain and out-of-domain (OOD) datasets. For the in-domain evaluation, we draw questions from the validation sets of VQAV2 and Oven, as well as the test set of ScienceQA. In addition, we generate unanswerable questions (UaVQA) as described in Section 3.2 and manually filter them for evaluation purposes. The final in-domain dataset consists of 1,000 samples.

For the OOD evaluation, we assess the model on three types of benchmarks: general VQA, knowledge-intensive VQA, and unanswerable VQA. For general VQA, we use the AOKVQA validation set (Schwenk et al., 2022), the GQA test set (Hudson and Manning, 2019), and the MM-Bench (en-dev) (Liu et al., 2023d). In the case of knowledge-intensive VQA, we employ the validation set of MMMU (Yue et al., 2024). For unanswerable VQA, we adopt the BeyondVisQA subset from MM-SAP (Wang et al., 2024b).

Baselines We consider both prompt-based and training-based methods. Refusal Prompt instruct the model to refuse answering when it lacks sufficient information by appending a prompt to the text input. The refusal prompt is: If you don't have enough information to answer the question, respond with "Sorry, I can not help with it." We also conduct supervised fine-tuning(SFT) as baseline. For questions within the extrinsic information boundary, the model is trained using the correct answers. For questions outside this boundary, since no correct answers exist, we assign an 'IDK' response as the label. By constructing the dataset in



Figure 4: Refusal rate and accuracy of models across different confidence levels. (a) Refusal Rate by Confidence: The model exhibits dynamic refusal behavior, with higher refusal rates for lower confidence levels and a tendency to answer directly for high-confidence questions. This indicates the model's awareness of its intrinsic information boundary. (b) Answered Accuracy by Confidence: The accuracy of the IDK-IT and CA-DPO models surpasses that of the original model, demonstrating that training methods focused on intrinsic boundary recognition improve the model's ability to provide accurate responses when choosing to answer.

7

this manner, we fine-tune models with about 11k instructions. Further details about evaluation can be found in the Appendix E.

Method	Acc	RefR	S_{trust}
LLaVA1.5-7B	12.00	46.10	-6.50
+Refusal Prompt	47.00	41.70	-4.10
+SFT	81.00	49.10	8.90
+IDK-IT	92.00	38.60	16.00
+CA-DPO	87.00	49.10	28.50
LLaVA1.5-13B	20.00	51.00	4.10
+Refusal Prompt	62.00	48.20	10.80
+SFT	78.00	53.60	17.30
+IDK-IT	89.00	41.80	21.20
+CA-DPO	93.00	49.00	32.00

Table 2: Performance on in-domain dataset. We present results on LLaVA1.5-7B and LLaVA1.5-13B. Bold values indicate the highest trustworthiness score.

5 Results and Analysis

5.1 Overall Results

The results on the in-domain datasets are presented in Table 2. Both IDK-IT and CA-DPO demonstrate notable improvements in trustworthiness scores compared to the baselines. Although IDK-IT does result in a decline in accuracy, it improves the model's trustworthiness by greatly increasing the model's refusal rate. In contrast, CA-DPO achieves a more balanced outcome by improving the refusal rate while maintaining model accuracy. This suggests that CA-DPO enables the model to better distinguish between known and unknown queries without sacrificing helpfulness. As illustrated in Table 1, IDK-IT and CA-DPO generalize well to OOD datasets. IDK-IT continues to mainly boost the refusal rate while CA-DPO strikes an effective balance between improving the refusal rate and preserving accuracy. 436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

5.2 The effectiveness of CA-DPO

To evaluate the effectiveness of CA-DPO, we used three types of preference data for the 'mixed' samples and trained LLaVA1.5 using the original DPO loss. These preference pairs include: (1) refusal > incorrect; (2) correct > incorrect; and (3) a combination of both. Table 3 presents the average performance on both in-domain and OOD datasets. The results indicate that models trained with the CA-DPO loss achieve a more balanced performance between accuracy and refusal rate, resulting in the highest trustworthiness score. This suggests that the CA-DPO method encourages the model to be more selective in its responses, striking an effective balance between providing helpful answers and refusing when necessary, thereby enhancing its overall trustworthiness.

5.3 Awareness of the Extrinsic Boundary

To comprehensively evaluate the model's awareness of the extrinsic information boundary, we conduct additional experiments on the unanswerable subset of VizWiz (Gurari et al., 2018) and the validation set of VQAv2-IDK (Cha et al., 2024).

421

422

423

494

425

426

427

428

429

430 431

432

433

434

Modal	Mathad	Data	I	n-Domai	in	Out-	Out-Of-Domain(Avg)			
WIOUEI	Method	Data	Acc	RefR	S_{trust}	Acc	RefR	S_{trust}		
LLaVA1.5-7B	DPO	(1)	47.50	30.20	25.20	50.30	29.46	30.06		
	DPO	(2)	51.00	26.30	28.30	55.08	18.62	28.78		
	DPO	(3)	49.50	26.30	25.30	52.88	20.35	26.11		
	CA-DPO	(3)	49.10	30.30	28.50	52.68	26.35	31.71		
LLaVA1.5-13B	DPO	(1)	47.10	36.10	30.30	52.71	24.14	29.56		
	DPO	(2)	49.60	31.40	30.60	56.37	16.45	29.20		
	DPO	(3)	48.60	33.90	31.10	55.19	20.83	31.21		
	CA-DPO	(3)	49.00	34.00	32.00	55.53	21.48	32.53		

Table 3: Performance comparison between models trained with different preference data using DPO and CA-DPO.

		Vizwiz(ua)	VQAv2-IDK(filter)
	original	9.00	2.80
LLaVA1.5-7b	IDK-IT	76.01	81.42
	CA-DPO	69.97	70.63
	original	9.60	2.60
LLaVA1.5-13b	IDK-IT	78.61	80.14
	CA-DPO	73.27	72.40

Table 4: Refusal Rate on unanswerable VQA datasets. VizWiz(ua) refers to the unanswerable subset of the VizWiz dataset, while VQAv2-IDK(filter) represents the filtered subset of VQAv2-IDK, where only questions with more than one 'IDK' annotation are retained.

VQAv2-IDK comprises questions from VQAv2 annotated with 'IDK' keywords. However, we observe that some of these questions can still be answered based on image information, suggesting that they remain within the extrinsic information boundary. Consequently, we filter out questions that contain only a single 'IDK' annotation.
VizWiz is a VQA dataset that includes visual questions posed by people who are blind. We select data labeled as 'unanswerable' to form its unanswerable subset. As shown in Table 4, our models appropriately provides refusal responses on these out-of-domain datasets, demonstrating their clear awareness of the extrinsic boundary.

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

5.4 Awareness of the Intrinsic Boundary

Although our overall results demonstrate that the 481 proposed training method significantly reduces mis-482 information by promoting refusals while maintain-483 ing performance, we aim to further investigate the 484 model's intrinsic awareness of its boundaries. To 485 486 this end, we analyze changes in the refusal rates relative to the confidence levels of LLaVA1.5-7B 487 on both the in-domain dataset and the AOKVQA 488 (OOD) dataset, as illustrated in Figure 4(a). The re-489 sults indicate that the model exhibits an awareness 490

of its own confidence, effectively refusing to answer when appropriate. For high-confidence questions, the model typically provides direct answers, while for lower-confidence questions, it demonstrates a higher likelihood of refusal. This adaptive refusal behavior reflects the model's capacity to distinguish between instances where it possesses sufficient knowledge and those where it does not, underscoring its intrinsic self-awareness. 491

492

493

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

525

526

Additionally, we calculate the accuracy of the answered questions, defined as: Answered Acc = $\frac{N_c}{N-N_r}$, where N_c is the number of correct answers and N_r is the number of refusals. Figure 4(b) shows the answered accuracy for the original model, the IDK-IT model, and the CA-DPO model. Notably, the accuracy curve for the original model is lower than those of both the IDK-IT and CA-DPO models. This indicates that training the model to recognize its intrinsic information boundaries through the IDK-IT and CA-DPO methods enhances its ability to more effectively utilize the information it possesses, leading to improved overall accuracy.

6 Conclusion

In this paper, we introduce the InBoL Framework to enhance the trustworthiness of MLLMs. By defining information boundaries, we create a data generation pipeline and apply novel training methods—IDK-IT and CA-DPO—to improve models' ability to avoid misinformation while maintaining helpfulness. Our user-centric evaluation approach also offers a simple way to assess trustworthiness. Experimental results show that our method effectively reduces misinformation and enhances model reliability, paving a feasible path for the future development of trustworthy MLLMs.

527

545

549

551

553

554

555

557

561

562

563

564

570

572

573

574

575

576

Limitations

In this work, we did not explore the generation 528 of explanations for refusal responses, an impor-529 tant and underexamined area. From the model's 530 perspective, many questions require reasoning pro-531 cesses to determine whether sufficient information is available to provide an accurate answer. By incorporating explanations for refusal responses, the 534 model could better learn when to refuse appropriately, thereby enhancing its awareness of its own limitations and boundaries. From the user's perspective, unexplained refusals may lead to confusion or dissatisfaction. Providing clear and in-539 terpretable justifications for refusals could make the refusal mechanism more transparent and user-541 friendly, significantly improving the overall user 542 experience. 543

> Actually, we found that our models not only can identify instances where sufficient information is lacking but also show potential in distinguishing between intrinsic and extrinsic information deficits. This suggests that it may be able to recognize the source of missing information and understand the reasoning behind refusal decisions. Detailed experimental results supporting this observation are presented in Appendix F.4.

Therefore, we plan to focus on enabling the model to generate well-reasoned and contextually appropriate refusal explanations for future work. This will involve developing methodologies for constructing relevant datasets and designing robust evaluation frameworks to assess the quality and relevance of the generated explanations. By making refusal responses more informative and transparent, we aim to further enhance the trustworthiness of the model while ensuring a more positive and engaging user experience.

References

- Alfonso Amayuelas, Kyle Wong, Liangming Pan, Wenhu Chen, and William Yang Wang. 2024. Knowledge of knowledge: Exploring known-unknowns uncertainty with large language models. In *Findings of the Association for Computational Linguistics ACL* 2024, pages 6416–6432, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Elmira Amirloo, Jean-Philippe Fauconnier, Christoph Roesmann, Christian Kerl, Rinu Boney, Yusu Qian, Zirui Wang, Afshin Dehghan, Yinfei Yang, Zhe Gan, et al. 2024. Understanding alignment in multi-

modal llms: A comprehensive study. *arXiv preprint arXiv:2407.02477*.

577

578

579

580

581

583

584

585

586

587

588

589

590

592

593

594

595

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

- 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.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*.
- Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou. 2024. Hallucination of multimodal large language models: A survey. *arXiv preprint arXiv:2404.18930*.
- Sungguk Cha, Jusung Lee, Younghyun Lee, and Cheoljong Yang. 2024. Visually dehallucinative instruction generation. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5510–5514. IEEE.
- Lida Chen, Zujie Liang, Xintao Wang, Jiaqing Liang, Yanghua Xiao, Feng Wei, Jinglei Chen, Zhenghong Hao, Bing Han, and Wei Wang. 2024. Teaching large language models to express knowledge boundary from their own signals. *arXiv preprint arXiv:2406.10881*.
- Qinyuan Cheng, Tianxiang Sun, Xiangyang Liu, Wenwei Zhang, Zhangyue Yin, Shimin Li, Linyang Li, Kai Chen, and Xipeng Qiu. 2024. Can ai assistants know what they don't know? *arXiv preprint arXiv:2401.13275*.
- Yunhao Fang, Ligeng Zhu, Yao Lu, Yan Wang, Pavlo Molchanov, Jang Hyun Cho, Marco Pavone, Song Han, and Hongxu Yin. 2024. vila²: Vila augmented vila. arXiv preprint arXiv:2407.17453.
- Chaoyou Fu, Haojia Lin, Zuwei Long, Yunhang Shen, Meng Zhao, Yifan Zhang, Xiong Wang, Di Yin, Long Ma, Xiawu Zheng, et al. 2024. Vita: Towards open-source interactive omni multimodal llm. *arXiv preprint arXiv:2408.05211*.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913.
- Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. 2018. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3608–3617.

Jiwoo Hong, Noah Lee, and James Thorne. 2024. Reference-free monolithic preference optimization with odds ratio. *arXiv preprint arXiv:2403.07691*.

631

632

634

641

647

651

652

660

671

672

673

674

675

677

678

679

681

- Hexiang Hu, Yi Luan, Yang Chen, Urvashi Khandelwal, Mandar Joshi, Kenton Lee, Kristina Toutanova, and Ming-Wei Chang. 2023. Open-domain visual entity recognition: Towards recognizing millions of wikipedia entities. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12065–12075.
- Drew A Hudson and Christopher D Manning. 2019. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024. Llavaonevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Yuxin Liang, Zhuoyang Song, Hao Wang, and Jiaxing Zhang. 2024. Learning to trust your feelings: Leveraging self-awareness in llms for hallucination mitigation. *arXiv preprint arXiv:2401.15449*.
- Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. 2023a. Aligning large multi-modal model with robust instruction tuning. *arXiv preprint arXiv:2306.14565*.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023b. Improved baselines with visual instruction tuning.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023c. Visual instruction tuning.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2023d. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS).*
- Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruti Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. 2024. Mm1: Methods, analysis & insights from multimodal llm pre-training. *arXiv preprint arXiv:2403.09611*.
- Chen Qian, Jie Zhang, Wei Yao, Dongrui Liu, Zhenfei Yin, Yu Qiao, Yong Liu, and Jing Shao. 2024. Towards tracing trustworthiness dynamics: Revisiting pre-training period of large language models. *arXiv preprint arXiv:2402.19465*.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36.

686

687

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

727

728

731

732

733

734

735

736

737

- Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022. A-okvqa: A benchmark for visual question answering using world knowledge. In *European conference on computer vision*, pages 146–162. Springer.
- Zhelun Shi, Zhipin Wang, Hongxing Fan, Zaibin Zhang, Lijun Li, Yongting Zhang, Zhenfei Yin, Lu Sheng, Yu Qiao, and Jing Shao. 2024. Assessment of multimodal large language models in alignment with human values. *arXiv preprint arXiv:2403.17830*.
- Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. 2023. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*.
- Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. 2024. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *arXiv preprint arXiv:2406.16860*.
- Fei Wang, Wenxuan Zhou, James Y Huang, Nan Xu, Sheng Zhang, Hoifung Poon, and Muhao Chen. 2024a. mdpo: Conditional preference optimization for multimodal large language models. arXiv preprint arXiv:2406.11839.
- Yuhao Wang, Yusheng Liao, Heyang Liu, Hongcheng Liu, Yu Wang, and Yanfeng Wang. 2024b. Mmsap: A comprehensive benchmark for assessing selfawareness of multimodal large language models in perception. *arXiv preprint arXiv:2401.07529*.
- Hongshen Xu, Zichen Zhu, Da Ma, Situo Zhang, Shuai Fan, Lu Chen, and Kai Yu. 2024. Rejection improves reliability: Training llms to refuse unknown questions using rl from knowledge feedback. *arXiv* preprint arXiv:2403.18349.
- Jing Xu, Andrew Lee, Sainbayar Sukhbaatar, and Jason Weston. 2023. Some things are more cringe than others: Preference optimization with the pairwise cringe loss. *arXiv preprint arXiv:2312.16682*.
- Yuqing Yang, Ethan Chern, Xipeng Qiu, Graham Neubig, and Pengfei Liu. 2023. Alignment for honesty. *arXiv preprint arXiv:2312.07000.*
- 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. *arXiv preprint arXiv:2408.04840*.

Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023. Do large language models know what they don't know? In Findings of the Association for Computational Linguistics: ACL 2023, pages 8653–8665, Toronto, Canada. Association for Computational Linguistics.

739

740

741

742

743

745

746

747

748

749

750

751

752

753

757

758

759

760

761

762

764

768

770

771

773

775

776

777

778

779

782

786

788

789

790

791

792

793

796

- Tianyu Yu, Yuan Yao, Haoye Zhang, Taiwen He, Yifeng Han, Ganqu Cui, Jinyi Hu, Zhiyuan Liu, Hai-Tao Zheng, Maosong Sun, et al. 2024a. Rlhf-v: Towards trustworthy mllms via behavior alignment from finegrained correctional human feedback. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13807–13816.
- Tianyu Yu, Haoye Zhang, Yuan Yao, Yunkai Dang, Da Chen, Xiaoman Lu, Ganqu Cui, Taiwen He, Zhiyuan Liu, Tat-Seng Chua, et al. 2024b. Rlaifv: Aligning mllms through open-source ai feedback for super gpt-4v trustworthiness. *arXiv preprint arXiv:2405.17220*.
- Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston.
 2024. Self-rewarding language models. *arXiv* preprint arXiv:2401.10020.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of CVPR*.
- Hanning Zhang, Shizhe Diao, Yong Lin, Yi Fung, Qing Lian, Xingyao Wang, Yangyi Chen, Heng Ji, and Tong Zhang. 2024a. R-tuning: Instructing large language models to say 'I don't know'. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 7113–7139, Mexico City, Mexico. Association for Computational Linguistics.
- Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong Duan, Bin Wang, Linke Ouyang, et al. 2024b. Internlmxcomposer-2.5: A versatile large vision language model supporting long-contextual input and output. *arXiv preprint arXiv:2407.03320*.
- Peng Zhang, Yash Goyal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2016. Yin and yang: Balancing and answering binary visual questions. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 5014–5022.
- Qinyu Zhao, Ming Xu, Kartik Gupta, Akshay Asthana, Liang Zheng, and Stephen Gould. 2025. The first to know: How token distributions reveal hidden knowledge in large vision-language models? In *European Conference on Computer Vision*, pages 127– 142. Springer.

Weihong Zhong, Xiaocheng Feng, Liang Zhao, Qiming Li, Lei Huang, Yuxuan Gu, Weitao Ma, Yuan Xu, and Bing Qin. 2024. Investigating and mitigating the multimodal hallucination snowballing in large vision-language models. In *Proceedings of the* 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 11991–12011, Bangkok, Thailand. Association for Computational Linguistics. 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

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

A Related work

A.1 MLLMs alignment

MLLM alignment seeks to reduce hallucinations and generate responses that are more closely aligned with human preferences through supervised fine-tuning and preference optimization. Tong et al. (2024); Li et al. (2024); Ye et al. (2024) enhance the perceptual and understanding capabilities of MLLMs by curating higher-quality visual instruction-tuning data. Fang et al. (2024) introduces a self-augmenting process that generates its own instructions to improve dataset quality. Reinforcement Learning and Direct Preference Optimization (Rafailov et al., 2024) have emerged as leading approaches for alignment, with recent advancements leveraging these methods to address visual hallucination issues. Sun et al. (2023) collects human preferences and adapts RLHF for multimodal alignment, while Yu et al. (2024a) improves MLLM performance by aligning model behavior through fine-grained human feedback corrections. Yu et al. (2024b) proposes a novel framework for gathering high-quality feedback data and uses an online feedback learning algorithm for model alignment. Additionally, Wang et al. (2024a) introduces a multimodal DPO objective that optimizes both image and language preferences, avoiding the overprioritization of language-only preferences.

A.2 Improving trustworthiness by Refusal

With the increasing capabilities of foundational models and the growing prevalence of AI agents, the trustworthiness of (multimodal) large language models has garnered significant attention. For LLMs, researchers primarily focus on the reliability of the model's knowledge, aiming for models to acknowledge their limitations and refuse to answer when encountering unknown knowledge. Yang et al. (2023) construct an honesty alignment dataset based on models' knowledge boundaries, replacing incorrect or uncertain LLM responses with "I don't know," and fine-tuning the model on

Refusal Template Sorry, I'm not sure about the answer to this question. Sorry, I can't provide a definite answer to this query. Sorry, I'm unable to provide an accurate response to this. Sorry. I can't help with that. Sorry, I can't help with this. Sorry, I can't help with it. Sorry, I can't answer that with the information I have. Sorry, I need more context to answer this question properly. Sorry, this is outside my scope of knowledge. Sorry, I'm not equipped to answer that question. Sorry, I don't have the data to respond to this guery. Sorry, I don't know the answer to that. Sorry, I don't have the necessary information to answer this. Sorry, I'm unable to provide a suitable answer. Sorry, I can't answer this question with certainty. Sorry, I don't have an answer for that right now. Sorry, I don't have enough information to answer that question. Sorry, I can not help with that. Sorry, I can not help with this. Sorry, I can not help with it.



878

this data. Cheng et al. (2024) proposed the concept of "Knowledge Quadrants," constructed the IDK dataset, and applied supervised fine-tuning (SFT) as well as preference-aware optimization to help models recognize their intrinsic knowledge boundaries. Zhang et al. (2024a) introduced Rtuning, which involves constructing and fine-tuning on a refusal-aware dataset, enhancing model's capabilities to refuse answering appropriately. Chen et al. (2024) directly judged whether the knowledge lies within the boundaries based on the model's intrinsic state and constructed training data to help the model express these boundaries. Liang et al. (2024) and Xu et al. (2024) employed Reinforcement Learning from Knowledge Feedback to teach models to refuse questions outside their knowledge boundaries, thus reducing hallucinations.

In multimodal scenario, only few works have considered the issue of refusal to answer. Unlike unimodal models, which focus on intrinsic boundaries, MLLMs mainly concentrate on the challenge of unanswerable questions. Liu et al. (2023a) proposed three types of negative instructions involving misleading or false premises in images, which models must learn to refuse. Cha et al. (2024) introduce the VQAv2-IDK dataset, which also annotates questions with "I don't know" answers to train models to appropriately refuse to respond when faced with unanswerable or ambiguous questions. Additionally, Shi et al. (2024) and Wang et al. (2024b) included subsets with unanswerable questions to evaluate the trustworthiness of MLLMs. Despite these advances, no prior work has systematically considered both the intrinsic boundaries of models and the extrinsic information provided in the input. Therefore, we propose the I-BaLF framework, which holistically integrates both aspects to guide MLLMs in refusing to answer when appropriate, thus significantly improving their overall trustworthiness. 879

880

881

882

883

884

885

887

888

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

B Refusal Template

Figure 5 shows the refusal template mentioned in Section 3.2.

C Detail of Data Construction

C.1 Hybrid Evaluator

We found that simple string matching—checking if the correct answers appear in the generated responses—was insufficient, as the model sometimes generates semantically similar answers that differ in wording. To address this issue, we propose to employ hybridized string matching and LLM-based evaluation methods to evaluate the accuracy of models. During our hybrid evaluation, we first use string matching to filter the model outputs that contain the exact ground truth answer and use Llama2-13B to check whether the remainder contains phrases that express the same semantic meaning.

To verify the effectiveness of our hybrid evaluator, we randomly sample 200 MLLM outputs from the our in-domain datasets for human annotators to

Prompt for Llama2-13B to evaluate the output result of MLLM

<s>[INST] <<SYS>>

Given a question and its ground truth answer, you should output True only if the statement includes the word or phrase of the ground truth answer. Otherwise, output False. Remember, only output True or False without any other words. <</SYS>>

Question: {question}

Ground truth answer: {answer}

MLLM Statement: {state}

Your judgment: [/INST]

908

909

910

911

912

913

914

915

916

917

918

920

921

Figure 6: LLM prompt for our hybrid evaluation, We use Llama2-13B for LLM evaluation.



Figure 7: Example 1 and Example 2 demonstrate the effectiveness of our hybrid evaluator. "England" and "the United Kingdom" represent the same country, and citrus and oranges denote the same fruit. Using string matching alone can not identify the correct answer from MLLM output, while our hybrid evaluator can effectively avoid false negatives.

assess the consistency between our hybrid evaluator and human evaluation. In these 200 samples, the Cohen's Kappa coefficient between our hybrid evaluator and human evaluator is 0.885, which is significantly higher than the coefficient of 0.749 for string matching. This result demonstrates the strong alignment between our method and human judgment.

The prompt used for our hybrid evaluation is shown in Figure 6.

The two cases shown in Figure 7 further demonstrate that our hybrid evaluator can correctly identify the sample with valid answers but ignores it when using only the string matching method.

C.2 Unanswerable Questions Generation

We define three criteria for determining whether a 923 question is unanswerable and instruct GPT-40 to verify if the questions meet these criteria. First, 925 questions may refer to subjects not present in the image, making it impossible to answer based on vi-927 sual information. Second, questions might include 929 incorrect premises about the subjects in the image, leading to misleading or unanswerable scenarios. Third, some questions may require additional con-931 text or information that cannot be inferred from the image alone. Figure 8 shows our prompt for gpt-40 933

to generate the unanswerable questions.

Based on these three reasons, we design corresponding questions to assess whether the generated questions are indeed unanswerable.

1. Does this question inquire about subjects that are not depicted in the image?

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

- 2. Does this question include an incorrect or misleading premise?
- 3. Does this question ask for information that is not available in the image?

Given the generated question and its corresponding image, we prompt GPT-4 to verify whether the question meets the specified criteria. Questions that receive a 'no' for all three criteria are filtered out. Additionally, we prompt the original model to generate a response to the unanswerable questions. If the model refuses to answer, those questions are also excluded from our dataset.

Figure 9 illustrates examples of unanswerable questions generated based on the proposed method. These examples demonstrate the diversity of scenarios leading to unanswerable questions, such as nonexistent objects or insufficient visual information. <Image> <TASK> Please generate two unanswerable questions based on the provided image. These questions should be related to the image but impossible to answer correctly based on the image. Specifically, create: 1. One question that refers to existing subjects with incorrect premises. 2. One question that requires information not derivable from the image. Additionally, for each question, provide a response beginning with 'Sorry' to explain why it cannot be answered based on the image. </TASK> <OUTPUT FORMAT> Question: {The generated question 1} Explanation: {The corresponding explanation 1}

Prompt for GPT-4o to Generate Unanswerable Questions

Question: {The generated question 2} Explanation: {The corresponding explanation 2}

</OUTPUT FORMAT>

Figure 8: Prompt for GPT-40 to Generate Unanswerable Questions

D Training Detail

958

961

962

963

964

965

966

967

969

970

972

973

974

975

976

978

982

985

For our experiments, we utilize the 7B and 13B versions of LLaVA-v1.5 as base models. We set $\delta_k = 0.8$ and $\delta_{uk} = 0.2$. The instruction dataset consists of 11k samples, with approximately 25% of the responses labeled as 'IDK.' Additionally, we generate around 24k preference pairs, with the ratio of unknown, mixed, and known samples approximately 1:1:2. For preference optimization, we first train the model on the IDK dataset and then conduct CA-DPO. LoRA is used for model training, with the LoRA rank r and α set to 16 and 32, respectively. The batch size is 16, and the learning rate is 2e-4, with training conducted for one epoch.

E Evaluation Detail

We here describe the used datasets:

- 1. **VQAv2** (Antol et al., 2015) is a widely-used dataset containing open-ended questions related to images, aimed at evaluating visual question answering.
- 2. **OVEN (Hu et al., 2023)** contains opendomain visual entity questions based on Wikipedia entries, requiring the model to possess extensive visual knowledge to provide accurate answers.
- 3. ScienceQA (Lu et al., 2022) comprises multimodal, multiple-choice questions across a diverse array of scientific topics.

4. **AOKVQA (Schwenk et al., 2022)** is a crowdsourced dataset featuring a wide range of questions that demand a broad understanding of commonsense and world knowledge.

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1001

1002

1004

1005

1007

1008

1009

1010

- 5. GQA (Hudson and Manning, 2019) is a dataset for real-world visual reasoning and compositional question answering. is a dataset designed for real-world visual reasoning and compositional question answering.
- 6. **MMMU** (Yue et al., 2024) is a benchmark developed to assess multimodal models across a variety of complex, multidisciplinary tasks that require college-level subject knowledge and advanced reasoning.
- 7. **MMBench (Liu et al., 2023d)** is a comprehensive benchmark for evaluating the multimodal capabilities of MLLMs, featuring questions that challenge both reasoning and perception.
- 8. **BeyondVisQA** (Wang et al., 2024b) is specifically designed to evaluate the self-awareness of MLLMs, particularly their ability to recognize "known unknowns." The questions in this dataset require information beyond the information provided by the input images.

To construct the in-domain evaluation dataset,1011we sample questions from the validation sets of1012VQAV2 and Oven, as well as the test set of Sci-1013enceQA. Importantly, we balance the confidence1014scores of these sampled questions to ensure that the1015



Figure 9: Examples of unanswerable questions.



Figure 10: Impact of confidence thresholds on performance. The heatmap displays the average trustworthiness scores across in-domain and OOD datasets, with scores normalized for comparison. The upper-right region, marked with a red box, demonstrates higher performance compared to other areas.

accuracy of LLaVA1.5-7B is approximately 50%. Additionally, we generate unanswerable questions (UaVQA) and manually filter them for use in the evaluation.

1016

1018

1019

1020

1021

1022

1024

1026

For both the MMMU and MMBench datasets, we use the following prompt for evaluation: "Answer with the letter corresponding to the correct option from the given choices." In contrast, for the remaining open-ended datasets, we presented only the questions, without any additional prompts. We use the proposed hybrid evaluator to assess the accucary for in-domain dataset, and we directly use the string matching strategy for the OOD dataset for simplicity.

F Supplementary experiments

F.1 Confidence threshold

We conducted experiments to analyze the impact of the confidence thresholds δ_k and δ_{uk} . Both thresholds were varied within the range $\delta_k, \delta_{uk} \in$ $\{2, 3, 4, 5, 6, 7, 8\}$, ensuring that $\delta_k > \delta_{uk}$. Using different combinations of these values, we generated an 'IDK' instruction dataset to fine-tune LLaVA1.5-7B. The results are illustrated in Figure 10, which displays average trustworthiness scores across in-domain and OOD datasets, with scores normalized for comparison.

We can see that the performance in the upperright region, highlighted by a red box, is notably higher than in other areas. Specifically, the combination of $\delta_k = 8$ and $\delta_{uk} = 2$ yields the best performance. This suggests that including data with intermediate confidence scores may not be beneficial for optimal model performance.

F.2 Data Composition

To achieve a balance between increasing the refusal rate and maintaining accuracy, we carefully adjust the proportions of "unknown," "mixed," and

1049

1050

1052



Figure 11: Data Composition Analysis. (a) For IDK-IT, varying the ratio of "known" data shows that a proportion of 0.75 yields the highest trustworthiness score across in-domain and out-of-domain datasets. (b) For CA-DPO, adjusting the ratio of "unknown," "mixed," and "known" data to 1:1:2 achieves the optimal balance between accuracy and refusal rate, as reflected in the trustworthiness score.

"known" data during training. All experiments in this section were conducted using the LLaVA1.5-7B model. For IDK-IT, we fixed the total training data size at 11K samples and varied the proportion of "known" data to balance accuracy and refusal rate. As shown in Figure 11(a), the trustworthiness score is highest when the proportion of "known" data is 0.75. This balance is consistent across both in-domain and out-of-domain (OOD) datasets. For CA-DPO, the data consists of three components: "unknown," "mixed," and "known." To simplify the experiment and focus on balancing accuracy and refusal rate, we fixed the ratio of "unknown" to "mixed" data at 1:1 and adjusted only the proportion of "known" data. As shown in Figure 11(b), the optimal trustworthiness score is achieved when the data ratio is 1:1:2, indicating a well-balanced trade-off between accuracy and refusal rate.

1053

1054

1055

1056

1058

1059

1060

1061

1063

1064

1065

1067

1070

1071

1072

1073

1075

1076

1079

1081

1083

F.3 Generalization of the Data Generation Pipeline

In Section 3.2, we introduced our data construction pipeline, which leverages a closed-source MLLM (GPT-40) to generate and filter unanswerable questions. A natural concern arises regarding the pipeline's reliance on GPT-40 and whether similar results can be achieved using other MLLMs, particularly open-source ones. To evaluate the generalizability of our pipeline, we employed an opensource MLLM (Qwen2-VL-72B) to generate unanswerable questions and used the resulting data to train LLaVA1.5-7B. The results shown in Table 5 and 6 demonstrate that the performance with data generated by Qwen2-VL-72B is comparable to that achieved with GPT-40. This finding suggests that our pipeline is flexible and can operate effectively with open-source MLLMs, making it more accessible and reproducible. 1084

1085

1086

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

F.4 Distinguishing Intrinsic and Extrinsic Information Deficits

As mentioned in Section 3.1, we categorize unknown questions into two types based on whether the model fails to answer due to a lack of visual information or internal knowledge. A key question is whether the model can differentiate between these two cases—that is, whether it can identify the specific type of information it lacks. To investigate this question, we apply linear probing to explore whether the model's internal representations encode features that differentiate between these two types of unknown cases.

Linear probing has been widely used in understanding and extracting knowledge of LLMs and MLLMs (Zhao et al., 2025; Qian et al., 2024). We selected 2,000 unknown questions from the training dataset with confidence score below 2, which are beyond extrinsic or intrinsic boundaries. For each question, we fed it into the MLLM and extracted the first generated token from the final layer of the model. We then trained a three-layer linear classifier to categorize these tokens into two classes. To evaluate the classifier, we used a set of 200 unknown questions (confidence scores below

Method Model for		A	OKVQ	A		GQA			MMMU			MMBench(en-dev)		
Wittillu	data generation	Acc	RefR	S_{trust}	Acc	RefR	S_{trust}	Acc	RefR	S_{trust}	Acc	RefR	S_{trust}	
IDK-IT	GPT-40	55.50	36.24	47.24	50.46	23.88	24.81	15.22	69.67	0.11	46.39	39.09	31.87	
IDK-IT	Qwen2-VL-72B	57.12	32.31	46.55	49.95	25.25	25.15	15.11	70.11	0.33	50.95	32.47	34.36	
CA-DPO	GPT-40	72.23	17.64	62.10	60.41	12.95	33.77	19.67	56.67	-4.00	58.42	18.13	34.97	
CA-DPO	Qwen2-VL-72B	71.79	20.35	63.93	58.32	15.25	31.89	21.11	50.33	-7.44	58.08	21.74	37.89	

Table 5: Performance on OOD datasets using GPT-40 and Qwen2-VL-72B for data generation.Bold values highlight the highest trustworthiness scores for each dataset.

	GPT-40-g	enerated data	Qwen2-VI	-generated data
	IDK-IT	CA-DPO	IDK-IT	CA-DPO
Vizwiz(ua)	76.01	69.97	74.49	71.39
VQAv2-IDK(filter)	81.42	70.63	79.25	75.22
BeyondVisQA	75.25	67.75	72.50	69.50

Table 6: Performance on unanswerable VQA datasets using GPT-40 and Qwen2-VL-72B for data generation.Bold values highlight the highest refusal rate for each dataset.

2) from the in-domain test set. The classification accuracy of the linear probing is shown in Table 7.

The results demonstrate that a simple linear classifier achieves high accuracy in distinguishing between intrinsic and extrinsic deficits based on the first generated token of MLLMs. This suggests that the model's internal representations encode features that effectively differentiate between these two types of unknown questions. This finding indicates that the model has the potential to recognize the source of missing information, even though its explicit refusal responses do not currently articulate this distinction. Incorporating this capability into the model's outputs could further enhance trustworthiness by providing more transparent refusals. Future work could leverage this potential by introducing mechanisms for explanation behind refusals, thereby aligning the model's behavior more closely with user expectations for trustworthy AI systems.

Model	Classification Acc(%)
LLaVA1.5-7B + IDK-IT	87.01
LLaVA1.5-7B + CA-DPO	80.83

Table 7: Classification accuracy for distinguishing between intrinsic and extrinsic deficits. Tokens were extracted using LLaVA1.5-7B + IDK-IT or LLaVA1.5-7B + CA-DPO.

1134 G Case Study

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127 1128

1129

1130

1131

1132 1133

1135Figure 12 presents examples comparing the per-1136formance of LLaVA1.5-7B with models trained1137using our framework. Models trained with IDK-IT

and CA-DPO are capable of refusing to answer1138questions where the original model makes errors.1139Notably, the CA-DPO model occasionally provides1140correct answers.These results demonstrate that1141our framework effectively mitigates the generation1142of misinformation, thereby enhancing the model's1143trustworthiness.1144

	Question: Which vehicle holds the most people? Ground Truth Answer: <u>Train</u> LLaVA1.5-7b: The vehicle holding the most people in the image is the blue car. +IDK-IT: Sorry, I don't know the answer. +CA-DPO: The <u>train</u> holds the most people in the image. It is a long passenger train with multiple cars, and it is filled with people.	
	Question: How many suitcases have tags? Ground Truth Answer: four LLaVA1.5-7b: There are three suitcases with tags in the image. +IDK-IT: Sorry, I don't know the answer to that. +CA-DPO: Sorry, I can't help with that.	
,	<	
	Question: How many suitcases have tags? A) Machine learning B) The highest mountain C) Teddy bear D) Sorry, I can't help with it. E) A little city Ground Truth Answer: D. Sorry, I can't help with it LLaVA1.5-7b: E +IDK-IT: Sorry, I can't help with it. +CA-DPO: Sorry, I can't help with it.	
/		
	Question: Who is depicted in «image 1> by Malcolm Drummond? A) Spencer Gore B) Harold Gilman C) Charles Ginne D) William Ratcliffe Ground Truth Answer: <u>C</u> LLaVA1.5-7b: D +IDK-IT: Sorry, I can't provide a definite answer to this query. +CA-DPO: Sorry, I can't answer this question with certainty.	
	Question: Who is depicted in «image 1> by Malcolm Drummond? A) Spencer Gore B) Harold Gilman C) Charles Ginne D) William Ratcliffe Ground Truth Answer: C LLaVA1.5-7b: D +IDK-IT: Sorry, I can't provide a definite answer to this query. +CA-DPO: Sorry, I can't answer this question with certainty.	
<image/> <image/>	Question: Who is depicted in «image 1> by Malcolm Drummond? A) Spencer Gore B) Harold Gilman C) Charles Ginne D) William Ratcliffe Ground Truth Answer: C LLaVA1.5-7b: D +IDK-IT: Sorry, I can't provide a definite answer to this query. +CA-DPO: Sorry, I can't answer this question with certainty. Question: Which of the following economic conditions was most important in creating the global trade network illustrated by the two diagrams? A) American demand for labor and African supply of slaves B) Japanese supply of silver and Middle Eastern demand for spices C) Western European demand for industrial raw materials and American and South Asian supply of cash crops D) Chinese and Indian demand for precious metals and European	
<image/>	Question: Who is depicted in «image 1> by Malcolm Drummond? A) Spencer Gore B) Harold Gilman C) Charles Ginne D) William Ratcliffe Ground Truth Answer: C LLaVA1.5-7b: D *IDK-IT: Sorry, I can't provide a definite answer to this query. *CA-DPO: Sorry, I can't answer this question with certainty. Question: Which of the following economic conditions was most important in creating the global trade network illustrated by the two diagrams? A) American demand for labor and African supply of slaves B) Japanese supply of silver and Middle Eastern demand for spices C) Western European demand for industrial raw materials and American and South Asian supply of cash crops D) Chinese and Indian demand for precious metals and European demand for Asian textiles and luxury goods.	

Figure 12: Examples illustrating the comparison between LLaVA1.5-7B and models trained with our framework (IDK-IT and CA-DPO)