# Do LVLMs Know What They Know? A Systematic Study of Knowledge Boundary Perception in LVLMs

Anonymous ACL submission

#### Abstract

001 Large Vision-Language Models (LVLMs) demonstrate strong visual question answering (VQA) capabilities but are shown to hal-004 lucinate. A reliable model should perceive its knowledge boundaries-knowing what it 006 knows and what it does not. This paper investigates LVLMs' perception of their knowledge 007 800 boundaries by evaluating three types of confidence signals: probabilistic confidence, answer consistency-based confidence, and verbal-011 ized confidence. Experiments on three LVLMs across three VQA datasets show that, although 012 LVLMs possess a reasonable perception level, there is substantial room for improvement. Among the three confidence, probabilistic and consistency-based signals are more reliable indicators, while verbalized confidence often 017 leads to overconfidence. To enhance LVLMs' perception, we adapt several established con-019 fidence calibration methods from Large Language Models (LLMs) and propose three effective methods. Additionally, we compare LVLMs with their LLM counterparts, finding that jointly processing visual and textual inputs decreases question-answering performance but reduces confidence, resulting in improved per-027 ception level compared to LLMs.

#### 1 Introduction

028

041

Large Vision-Language Models (LVLMs) are capable of processing both textual and visual information simultaneously, demonstrating strong performance on visual question-answering (VQA) task (Bai et al., 2025a; Wu et al., 2024; OpenAI et al., 2024). However, when faced with questions beyond their knowledge boundaries, LVLMs often hallucinate—generating seemingly plausible but factually incorrect responses (Liu et al., 2024a; Bai et al., 2025b). This is unacceptable in safetycritical domains such as healthcare. Knowing when an LVLM can answer correctly not only helps us determine when to trust the model but also enables adaptive retrieval-augmented generation (RAG), triggering RAG only when the model does not know the answer, which improves both the efficiency and effectiveness of RAG (Ni et al., 2024). 043

044

045

046

047

049

051

054

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

A trustworthy model should have a clear perception of its knowledge boundaries—knowing what it knows and what it does not. While this ability has been extensively studied in Large Language Models (LLMs) (Xiong et al., 2024; Tian et al., 2023; Moskvoretskii et al., 2025), it remains underexplored in LVLMs. A model's perception level is assessed by the alignment between its confidence and actual performance, with correctness of the answer serving as a proxy for performance. Therefore, the emphasis is on whether LVLMs can provide confidence that matches their performance. We focus on binary confidence because it directly helps us decide whether to trust the model.

In this work, we explore this question by examining three representative types of confidence signals that are widely used in LLMs: 1) Probabilistic confidence (Desai and Durrett, 2020; Guo et al., 2017a). The confidence is measured by the generation probability of tokens in the output. 2) Answer consistency-based confidence (Zhang et al., 2024; Manakul et al., 2023b). Some studies argue that token-level probabilities poorly reflect a model's semantic confidence and are not suitable for blackbox models. Instead, they suggest using semantic consistency across multiple responses as a confidence indicator. 3) Verbalized confidence (Lin et al., 2022; Yang et al., 2024b). The natural language confidence expressed by the model, offering an intuitive and model-agnostic signal without requiring repeated sampling.

We conduct experiments using three representative models—-Qwen2.5-VL (Bai et al., 2025a), DeepSeek-VL2 (Wu et al., 2024), and LLaVA-v1.5 (Liu et al., 2024b)—on three datasets: Dyn-VQA (Li et al., 2025b), MMMU Pro (Yue et al., 2024), and Visual7W (Zhu et al., 2016). Results show that LVLMs are able to perceive their knowledge bound-

110 111

- 112

113 114

115 116

117

119

118

120 121

123

125

122

124

128 129

126

127

130

131

132

133 134

**Related Work** 2

aries to some extent, but there remains considerable

room for improvement. Among the three types of

confidence, probabilistic and answer consistency-

based confidences are more aligned with LVLMs'

performance but rely on in-domain data for bina-

rization, while verbalized confidence have weaker

To enhance LVLMs' perception capabilities, we

adopt several representative confidence calibration

methods originally designed for LLMs and propose

three new approaches tailored for LVLMs: Img-

CoT, Prob-Thr, and Cross Model. Our results show

that methods which engage the model's reasoning

abilities can enhance both answer accuracy and

verbalized confidence perception-level, whereas

existing consistency-based methods have limited

effectiveness and do not generalize well to LVLMs.

Our proposed three methods are effective on differ-

Compared to LLMs, LVLMs need to process an

additional visual modality and integrate informa-

tion across different modalities. This raises a ques-

tion: how does the perception ability of LVLMs dif-

fer from that of LLMs? To investigate this, We com-

pare LVLMs with their corresponding LLMs on

the Dyn-VQA dataset (Li et al., 2025b; Tian et al.,

2023). This dataset provides parallel visual-textual

and pure textual queries, ensuring fair comparison

between LLMs and LVLMs. We focus on verbal-

ized confidence because it can reflect the model's

language capabilities. Experimental results show

that: 1) LVLMs exhibit lower VQA performance

but higher perception accuracy compared to their

LLM counterparts. 2) Certain prompting methods

are ineffective for LVLMs, showing that LVLMs

by the following two reasons: 1) Compared to pro-

cessing single-modality information, handling vi-

sual information and integrating two modalities

pose challenges for LVLMs, which leads to low

VQA performance. However, this also leads to

more conservative confidence expression, resulting

in lower overconfidence and a more accurate per-

ception level. 2) Training LLMs without enough ca-

pacity to handle additional visual information may

lead to a loss of their language abilities, thereby im-

pairing their instruction following ability. Through

controlled experiments on different model scales

and different input modalities, our results supported

these hypotheses.

We hypothesize these phenomena may be caused

have weaker instruction-following capabilities.

ent datasets and models.

alignment and tends to be overconfident.

LLM Knowledge Boundary Perception. Prior research has primarily focused on knowledge boundary perception in LLMs, with various methodologies proposed to elicit confidence: verbalized confidence, where models directly articulate their confidence (Yang et al., 2024b; Yin et al., 2023; Zhang et al., 2023); consistency based confidence that derive confidence from answer consistency across multiple samples (Manakul et al., 2023a; Agrawal et al., 2024); probabilistic confidence, leveraging generated token likelihoods (Guo et al., 2017b; Ma et al., 2025); and internal state probing confidence, examining hidden states (Azaria and Mitchell, 2023; Ni et al., 2025). Differently, our work investigates knowledge boundary perception in LVLMs and provides the first systematic comparison of these methods in the multimodal setting.

LVLMs. Previous studies have established the widespread adoption of LVLMs in safety-critical domains such as healthcare (Li et al., 2023; Hu et al., 2024) and autonomous driving (Cui et al., 2024; Jiang et al., 2024). While these applications demonstrates LVLMs' functional capabilities, studies show LVLMs frequently produce hallucinations (Bai et al., 2025b; Sahoo et al., 2024). The current body of work investigates this limitation on different aspects. Some work surveys hallucination types and their causes (Liu et al., 2024a; Zhou et al., 2024; Lan et al., 2024), while others focus on mitigating hallucinations (Li et al., 2025a; Wang et al., 2024a; Xiao et al., 2025). A distinct but less explored research thread investigates LVLMs' knowledge boundary as a potential framework for enhancing model reliability (Chen et al., 2025; Wang et al., 2024b; Leng et al., 2024). We take this line of work a step further by introducing a novel comparative paradigm that compares perception between LVLMs with their LLM counterparts.

#### 3 **Preliminaries**

In this section, we provide an overview of our task.

#### 3.1 **Task Formulation**

Visual Question Answering. The goal of visual question answering (VQA) can be described as follows. For a given question q and an image i, the model is asked to provide an answer a based on the question q and image i, that is,  $a = f_{model}(q, i)$ .

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

177

178

179

180

181

**LVLM Knowledge Boundary Perception.** We assess the perception of LVLM's knowledge boundary with the alignment between confidence and its actual performance. Here, we use the model's visual question answering correctness to serve as a proxy for performance, and elicit different kinds of model confidence estimates.

182

183

184

188

189

190

191

192

193

195

196

197

198

199

201

203

204

207

208

210

211

212

214

215

216

217

218

219

221

222

224

**Confidence Estimation.** In this paper, we conduct experiments on the following three kinds of model confidence estimates. As widely adopted and training-free, they can be elicited without changing the internal knowledge of models.

Probabilistic confidence is elicited through the aggregation of token probabilities for scoring, followed by applying a threshold to binarize the score into confidence. It is efficient but only captures lexical-level confidence and requires threshold tuning on a held-out set, which leads to poor generalizability. Some studies also argue that it is not applicable to black-box models (Kuhn et al., 2023).

Answer consistency-based confidence is elicited by calculating the consistency of multiple generated responses. The core idea is that if the model knows the correct answer, multiple sampled answers should be semantically consistent. It better captures semantics than probabilistic confidence, but is computationally expensive and still requires fitting a threshold (Manakul et al., 2023a).

Verbalized confidence is elicited by directly asking the model to express confidence (Yang et al., 2024b). Compared to the other two confidences, this confidence reflects models' self-awareness of their knowledge boundaries. Moreover, it eliminates the need for threshold fitting and multiple sampling. Therefore, this kind of confidence receives our primary focus.

# 4 Knowledge Boundary Perception in LVLMs

This section introduces experimental setup to evaluate LVLMs knowledge boundary perception ability. Along with the elicited confidence and confidence calibration methods evaluated by us.

### 4.1 Existing Methods

Here, we systematically introduce three basic confidence estimates in Section 3, along with several confidence calibration methods originally designed for LLMs. We also propose new methods. Detailed prompts are in Appendix A.1. Basic confidence estimates are in <u>underline</u>, and others are existing confidence calibration methods.

## 4.1.1 Vanilla Confidence Estimation Methods

231

233

234

235

236

237

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

Probabilistic confidence is elicited through token probabilities. Here, we focus on the output perplexity of models.

• **Perplexity Threshold (PPL-Thr)**: Perplexity quantifies a model's uncertainty in content generation (Cooper and Scholak, 2024). We binarize this metric into confidence by applying a threshold decided on a held-out set.

Answer consistency-based confidence requires models to generate multiple responses, compute their consistency, and apply a threshold to the consistency scores for confidence elicitation.

• **Random Sample (Random)**: Simply sample responses without modifying input.

We evaluate two types of verbalized confidence: (1) Single-step verbalized confidence, which is generated simultaneously with the answer, and (2) Double-step verbalized confidence, which is generated by asking the model for an answer in the initial round of dialogue, then providing its confidence in the second round. The distinction between them lies in cognitive focus allocation: single-step confidence elicitation demands concurrent attention to both answer and confidence generation, whereas double-step confidence elicitation enables sequential processing.

- Single-step Vanilla (Vanilla) : Simply ask the model to generate both the answer and confidence in a single interaction.
- **Double-step Self Judging (Self-Jud)**: First, acquiring the model to provide an answer to the question, then asking it to generate confidence.

### 4.1.2 Calibrating Verbalized Confidence

The four methods below aim to calibrate singlestep verbalized confidence:

- Chain-of-Thought (CoT): Zero-shot Chain-of-Thought prompting, Applying "Analyze step by step" to the query (Kojima et al., 2023).
- **Punish**: Penalizing overconfidence via the instruction "You will be punished if the answer is not right but you say certain".
- **Explain**: Requesting models to provide answer explanations before generating their confidence. The four methods below aim to calibrate double-step verbalized confidence:
- Chain of Thought (CoT): Applying the Chainof-Thought prompt in the confidence elicitation round of dialogue.

- 329 330 331
- 332
- 334

336

337

338

339

340

341

342

343

346

347

348

349

350

351

353

354

355

357

358

359

360

361

362

364

365

366

368

• **Challenge**: We prepend the critical prompt "I don't think your answer is right" to the query in the confidence elicitation round in order to guide the model to be less overconfident.

281

287

290

291

297

298

305

306

307

311

313

314

315

316

• **Punish**: Applying the Punish prompt in the confidence elicitation round of dialogue.

#### 4.1.3 Calibrating Answer Consistency-Based Confidence

- **Rephrasing (Rephr)**: To address persistent errors caused by a specific question phrase, rephrase the original question into semantically equivalent variants with different phrases (Yang et al., 2024a).
- Noised Image (Noised-Img): Reducing persistent errors caused by a specific image by creating semantically equivalent variants through the addition of subtle noise to the original image.
  - **Rephrasing and Noised Image (Reph+Nois)**: A combination of the Rephrasing and the Noised Image methods.

#### 4.2 Newly Proposed Methods

- Image Chain of Thought (Img-CoT): Prompting models to generate textual image descriptions before reasoning in order to convert visual modality information to textual modality.
- **Probability Threshold (Prob-Thr)**: Prompting models to generate continuous probabilities of answers (0–1), then applies a threshold to them to generate binary confidences. The threshold is decided on a held-out set.
- **Cross Model**: Utilizing generated responses from different models to calculate the consistency score, this method can be viewed as using other models' answers to evaluate whether the answer generated by a given model is reliable.

### 4.3 Experimental Setup

Datasets. We conduct experiments on three VQA 317 benchmark datasets. They emphasize on LVLM's 318 different abilities. Visual7W (Zhu et al., 2016) emphasizes abilities in vision comprehension, it con-320 tains 70K image-QA pairs for basic visual under-321 standing. Dyn-VQA (Li et al., 2025b) emphasizes 322 language reasoning, it contains 1.5K questions testing multi-modal knowledge and multi-hop reason-324 ing.; MMMU Pro (Yue et al., 2024) emphasizes 325 both vision and language capability, it contains 12K 326 expert-curated multimodal questions. For evaluation, we respectively sample 550 questions from 328

Dyn-VQA and MMMU Pro datasets, and sample 500 questions from the Visual7W dataset.

**Models.** We conduct experiments on three representative LVLMs: Qwen2.5-VL-7B (Bai et al., 2025a), DeepSeek-VL2-16B (Wu et al., 2024), and LLaVA-v1.5-7B (Liu et al., 2024b).



Figure 1: Count of samples for various matches between answer correctness and model confidence. We use Total = FN + FP + TN + TP to represent the total number of samples.

**Metrics.** We mainly utilize the evaluation metrics proposed by Ni et al. (2024): (1) **Uncertain-Rate (Unc-R.)**:  $\frac{FN+TN}{Total}$  represents the proportion where the judgement of the answer is unconfident. (2) **Accuracy (Acc.)**:  $\frac{TP+FN}{Total}$  indicates the ratio of correct answers generated by the model. (3) **Alignment (Align.)**:  $\frac{TP+TN}{Total}$  represents the proportion of samples where confidence matches the result, we mainly use this metric to assess the model's knowledge boundary perception ability. (4) **Overconfidence (Overco.)**:  $\frac{FP}{Total}$  is the ratio of model-generated answer is incorrect, but the judgement is confident. (5) **Conservativeness (Conser.)**:  $\frac{FN}{Total}$  is the ratio of model-generated answer is correct but the judgement is unconfident.

#### 4.4 Results and Analysis

Table 1 presents the results of alignment performance across different datasets and models. Please refer to Appendix A.2 for implementation details and detailed results.

# 4.4.1 Performance of Different Types of Confidence

Here, we analyze three basic elicited confidence's performance. Our findings are as follows:

1) **Compared to verbalized and probabilistic confidence, answer consistency-based confidence often shows higher alignment.** As shown in Table 1, the basic answer-consistency based confidence (Random) achieves higher alignment compared to verbalized (Vanilla, Self-Jud) and probabilistic confidences (PPL Thr) on both LLaVA-1.5 and Deepseek-VL2. This may be because, unlike probabilistic confidence that operates at the lexical level, answer consistency-based confidence better

		Qwen2.5-V	L		LLaVA-1.5	;	DeepSeek-VL2				
method	Dyn-VQA	Visual7W	MMMU Pro	Dyn-VQA	Visual7W	MMMU Pro	Dyn-VQA	Visual7W	MMMU Pro		
<u>Vanilla</u>	0.7623	0.5840	0.4909	0.5338	0.4140	0.2509	0.6527	0.2820	0.2727		
СоТ	0.7824	0.6080	0.6818	0.5375	0.3940	0.2418	0.6362	0.5540	0.3836		
Punish	0.7112	0.5520	0.5000	0.4899	0.4180	0.3745	0.7093	0.3500	0.3145		
Explain	0.8117	0.6180	0.5782	0.4534	0.3900	0.2109	0.6984	0.5700	0.3491		
Img-CoT	0.7276	0.6060	0.7182	0.5484	0.4140	0.2964	0.6344	0.5360	0.5236		
Self-Jud	0.3272	0.5500	0.5609	0.2468	0.4220	0.3327	0.1993	0.4780	0.4236		
СоТ	0.6435	0.5700	0.5255	0.1463	0.4200	0.3218	0.2029	0.4760	0.4255		
Challenge	0.8080	0.5280	0.4891	<del>0.8995</del>	0.5800	<del>0.6782</del>	0.8007	0.5240	0.5709		
Punish	0.3272	0.5300	0.5164	0.1298	0.4200	0.3218	0.4936	0.5300	0.4345		
Prob-Thr	0.5960	0.5820	0.5855	0.7971	0.6140	0.6091	<u>0.6910</u>	0.6060	<u>0.5218</u>		
Random	0.5448	0.5700	0.5327	<u>0.8976</u>	0.7080	0.6709	0.8026	0.6460	0.6000		
Noised Img	0.7313	0.6000	0.5400	0.8958	0.6740	0.6655	0.8062	0.6300	<u>0.5818</u>		
Rephr	0.8026	0.5660	0.5364	0.8976	0.6920	0.6672	<u>0.8080</u>	0.6260	0.5764		
Reph+Nois	0.7733	0.5500	0.5509	0.9013	0.6780	0.6655	0.8099	0.6120	0.5618		
Cross Model	0.8208	0.6320	0.5800	<u>0.8976</u>	0.6520	0.6618	0.8062	0.6740	0.5964		
PPL Thr	0.7916	0.6020	0.6073	0.8519	0.7060	0.6800	0.7934	0.6280	0.5345		

Table 1: Performance of alignment on three datasets and three LVLMs. Best results of each kind of confidence in **bold** and second best in <u>underline</u>. Experimental observations show that LLaVA demonstrates a pattern of complete answer denial when being challenged. We therefore omitted these data from our results.

captures semantics by evaluating answer consistency (Kuhn et al., 2023), achieving higher alignment. Additionally, while verbalized confidence is uncalibrated, eliciting answer consistency-based confidence calibrating a threshold on a held-out set, further improves alignment.

Despite answer consistency-based confidence exhibiting high alignment, it comes at a cost: eliciting this kind of confidence requires generating multiple responses, incurring high computational costs. And its reliance on a held-out set for threshold calibration limits its generalizability.

2) **Probabilistic confidence surpasses verbal**ized confidence in alignment performance. As shown in Table 1, probabilistic confidence's alignment performance consistently surpasses verbalized confidence, and it outperforms answer consistency-based confidence on Qwen2.5-VL. Though it falls behind consistency-based confidence on LLAVA-1.5 and DeepSeek-VL2, the alignment differences are small. Additionally, it functions more efficiently without the high computational cost of generating multiple responses.

However, probabilistic confidence, like answer consistency-based confidence, still requires threshold calibration on a held-out set, which affects its generalizability.

3) Verbalized confidence demonstrates lower alignment compared to probabilistic and answer consistency-based confidences, and judges an-

	Dyn-	VQA	Visu	al7W	MMMU Pro				
method	Conser.	Overco.	Conser.	Overco.	Conser.	Overco.			
Vanilla	0.1024	0.1353	0.0900	0.3260	0.1327	0.3764			
СоТ	0.0786	0.1389	0.1340	0.2580	0.1127	0.2055			
Punish	0.0804	0.2084	0.0820	0.3660	0.1127	0.3873			
Self-Jud	0.0018	0.6709	0.0180	0.4320	0.0701	0.3692			
СоТ	0.0329 0.3236		0.0280	0.4020	0.1000	0.3745			
Punish	0.0018	0.6709	0.0120	0.4580	0.0200	0.4636			

Table 2: The performance of verbalized confidence on Qwen2.5-VL, single-step confidences are in blue and souble-step confidences are in orange.

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

**swers overconfidently.** Compared to probabilistic and answer consistency-based confidences, eliciting verbalized confidence is computationally efficient and generalizable. However, as shown in Table 1, both single-step (Vanilla) and double-step (Self-Jud) verbalized confidences' alignment are lower than the other two confidences. To investigate the cause of it, we calculate the conservativeness and overconfidence on verbalized confidence, as shown in Table 2, we find that the ratio of overconfident responses is substantially higher than conservative responses. This pattern suggests that LVLMs, like LLMs, are intrinsically biased toward affirming their own output (Groot and Valdenegro-Toro, 2024; Sun et al., 2025).

Table 2 also shows that double-step verbalized confidence exhibits more severe overconfidence than its single-step counterpart. This may be because the model's self-generated answers in the first round of dialogue serve as false positive signals of

394

		Dyn-VQ	4		Visual7V	V	MMMU Pro				
method	Acc.	Align.	Overco.	Acc.	Align.	Overco.	Acc.	Align.	Overco.		
Vanilla	0.1846	0.7623	0.1353	0.4380	0.5840	0.3260	0.4564	0.4909	0.3764		
СоТ	0.2121	0.7824	<u>0.1389</u>	0.4920	0.6080	0.2580	0.6436	0.6818	0.2055		
Img-CoT	0.2048	0.7276	0.2066	0.5020	0.6060	0.3080	0.6636	0.7182	0.1691		
Explain	0.1956	0.8117	0.0823	0.4740	0.6180	0.2720	0.5309	0.5782	0.2982		

Table 3: The performance of single-step reasoning elicitation methods on Qwen2.5-VL.

its capability, reinforcing overconfident behavior through misleading model to self-affirmation.

419

420

421

422

423

494

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

455

456

457

458

#### 4.4.2 Confidence Calibration in LVLMs

In this section, we evaluate the effectiveness of existing confidence calibration methods developed for LLMs in the context of LVLMs, as well as our proposed methods.

For existing confidence calibration methods, our observations are as follows:

1) Single-step reasoning elicitation methods effectively enhance the accuracy and alignment of LVLMs. As shown in Table 1, we found reasoning elicitation methods (Explain, CoT, and Img-CoT) exhibit high alignment. To further investigate them, we calculate other metrics about them. Table 3 shows that different reasoning elicitation methods excel on specific datasets: CoT method improves alignment and accuracy across all datasets and causes overconfidence on Dyn-VQA. The Explain method outperforms CoT in alignment on Visual7W and Dyn-VQA datasets. This observed difference may stem from the Explain method's design: while the CoT method enforces step-by-step reasoning, the Explain method prioritizes direct justification, thus reducing redundant context for simple questions and improving the calibration of LVLMs' confidence outputs.

2) Answer consistency-based confidence calibration methods improve alignment on Qwen2.5-VL, but show limited effectiveness on other models. We observed that, even when sampling responses at the same temperature of 1.0, models differ in their output diversity. As shown in Table 1, when random sampling Qwen2.5-VL's responses, it tends to generate consistent yet incorrect responses, resulting in low alignment. However, both the rephrasing and the noised image methods show effectiveness in mitigating this tendency, consequently achieving higher alignment. In contrast, LLaVA-1.5 and DeepSeek-VL2 generate more di-

verse outputs when the response is incorrect, allowing the Random Sampling method to perform well and making Noised Image and Rephrasing methods less effective in enhancing alignment by comparison. 459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

497

498

499

We propose Image Chain of Thought, Probability Threshold, Cross Model Consistency methods in Section 4.1, their performances are as follows:

1) Image Chain of Thought method effectively enhances alignment and accuracy on MMMU Pro. As shown in Table 3, Img-CoT demonstrates remarkable performance on the MMMU Pro dataset, which requires both strong visual perception and reasoning capabilities. It improves accuracy and alignment, outperforming CoT method. This indicates its converting visual modality into language modality mechanism can effectively enhance models' comprehension of the content in the image, thereby achieving this superior performance.

2) **Probability Threshold method shows higher alignment than other double-step verbalizated confidence calibration methods.** As shown in Table 1, the Probability Threshold method outperforms alternative double-step methods. Despite the need to calibrate the threshold, it effectively enhances alignment.

3) Cross Model method outperforms other answer consistency based methods on LLaVA-1.5 and Qwen2.5-VL. Table 1 demonstrates that the Cross Model method performs superior alignment on both LLaVA-1.5 and Qwen2.5-VL compared to other answer consistency-based confidence calibration methods. This improved performance validates the method's effectiveness, as it establishes a robust mechanism for evaluating answer correctness by leveraging responses from diverse models.

## 5 Perception Comparison Between LVLMs and LLMs

Compared to LLMs, LVLMs need to process additional visual modality and integrate information

			Qwen2.	5			VL2		LLaVA-1.5							
method	Model	Unc-R.	Acc	Align.	Conser.	Overco.	Unc-R.	Acc	Align.	Conser.	Overco.	Unc-R.	Acc	Align.	Conser.	Overco.
Vanilla	LVLM	0.782	0.185	<b>0.762</b>	0.102	<b>0.135</b>	<b>0.788</b>	0.146	<b>0.653</b>	<b>0.141</b>	0.207	<b>0.490</b>	0.088	<b>0.534</b>	<b>0.022</b>	0.444
	LLM	<b>0.788</b>	<b>0.285</b>	0.729	0.172	0.099	0.161	0.225	0.338	0.024	<b>0.638</b>	0.011	<b>0.141</b>	0.152	0.001	<b>0.848</b>
СоТ	LVLM	0.728	0.212	<b>0.782</b>	<b>0.079</b>	0.139	<b>0.638</b>	0.170	<b>0.636</b>	0.086	0.278	<b>0.512</b>	0.084	<b>0.538</b>	<b>0.031</b>	0.431
	LLM	0.448	0.294	0.651	0.046	<b>0.304</b>	0.095	0.296	0.362	0.015	0.623	0.117	<b>0.199</b>	0.302	0.007	<b>0.691</b>
Punish	LVLM	0.711	0.168	<b>0.711</b>	0.080	<b>0.208</b>	<b>0.848</b>	0.161	<b>0.709</b>	<b>0.150</b>	0.141	<b>0.450</b>	0.095	<b>0.490</b>	<b>0.027</b>	0.483
	LLM	<b>0.956</b>	0.294	0.713	<b>0.269</b>	0.018	0.266	0.229	0.455	0.020	0.525	0.057	<b>0.152</b>	0.201	0.004	0.795
Explain	LVLM LLM	<b>0.828</b> 0.536	0.196 <b>0.298</b>	<b>0.812</b> 0.673	<b>0.106</b> 0.080	0.082 0.247	<b>0.786</b> 0.079	0.168 0.252	<b>0.698</b> 0.320	<b>0.128</b> 0.006	0.174 <b>0.675</b>	<b>0.421</b> 0.159	0.084 <b>0.219</b>	<b>0.453</b> 0.364	<b>0.027</b> 0.007	0.519 <b>0.629</b>

Table 4: LLMs and LVLMs comparison for single-step verbalization based methods on Dyn-VQA.



Figure 2: Comparative analysis of instruction following ability across model scales.

across different modalities. This raises a question: how does the perception of LVLMs differ from that of LLMs? Knowing these distinctions is valuable for developing trustworthy LVLMs.

In this section, we investigate the difference of knowledge boundary perception between LVLMs and their LLM counterparts. Focusing on verbalized confidence cause it directly reflects models' self-awareness of their knowledge boundaries. We further propose several hypotheses about these differences' underlying causes and validate them through the comparison between different model scales and input modalities.

#### 5.1 Experimantal Setup

500

503

504

505

507

508

509

510

511

512

513

514

515

516

517

518

519

522

527

**Datasets.** In this section, we mainly focus on Dyn-VQA dataset. Dyn-VQA provides both VQA question image pairs and their semantically equivalent QA questions (e.g., QA: "How many humans have landed on Mars?" vs. VQA: "How many humans have landed on this planet?" with an image of Mars). This enables fair model performance comparison across text-only modality and vision-text modality inputs.

Models. In this section, we compare LVLMs with
their base LLM counterparts to ensure fair comparison: Qwen2.5-VL, DeepSeek-VL2, LLaVA-v1.5
vs Qwen2.5, DeepSeek-MoE, Vicuna-v1.5.

#### 5.2 Results and Analysis

Here, we apply VQA queries on LVLMs, and their semantic equivalent QA queries on LLMs to fairly

compare them. And focus on single-step verbalized confidence. We defer results about other kinds of confidence to Appendix A.3. Here are our findings: 530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

1) Compared to LLMs, LVLMs struggle to follow certain methods' instructions, leading to performance deviating from expected. As shown in Table 4, Qwen2.5-VL cannot effectively follow the Punish instruction. As a result, this method not only fails to reduce overconfidence but actually exacerbates it, leading to lower alignment than Vanilla. Similarly, LLaVA-1.5 disregards CoT and Explain instructions, persistently generating responses without proper reasoning or explanation, which results in lower accuracy. This stands in contrast to LLMs, where the Punish method effectively reduces Qwen2.5's overconfidence; CoT and Explain instructions reliably ignite reasoning responses in Vicuna-1.5, thus improving its accuracy.

2) For single-step verbalized confidence, LVL-Ms tend to have lower accuracy compared to LLMs. Along with higher alignment due to reduced overconfidence. As shown in Table 4, under all single-step verbalized confidence for the three series of models, the answer accuracy of LVLMs is lower than that of LLMs. Meanwhile, LVLMs exhibit a higher uncertain-rate compared to LLMs. Specifically, LLaVA exhibits an average accuracy reduction of 0.09 with a concurrent 0.382 increase in uncertain-rate than its counterpart LLM. And in DeepSeek-VL2, we observe an 0.089 accuracy decrement paired with a 0.615 surge in uncertainty than LLM. Compared to LLMs, LVLMs' accuracy drop is relatively smaller than their uncertain-rate increase, thus they demonstrate less severe overconfidence than LLMs, leading to relatively higher alignment in their responses.

561

562

563

566

569

571

572

573

574

578

580

581

583

585

586

588

589

591

596

599

605

607

610

#### 5.3 Analysis Across Model Scales and Modalities

Building upon the findings discussed in the previous subsection, we observe notable performance distinctions between LLMs and LVLMs, which motivate us to propose the following hypothesis regarding their potential underlying causes:

1) Model capacity bottleneck: We hypothesize that the inferior instruction-following abilities of LVLMs stems from their internal capacity limitations, where visual modality integration competes for models' internal parameter resources that would otherwise support language processing capabilities.

2) Cross-modal limitation awareness: While the LVLMs demonstrate lower accuracy than LLMs, their verbalized confidence shows better alignment with performance. We hypothesize this stems from two factors: (1) LVLMs' constrained cross-modal processing ability leads to degraded multimodal VQA accuracy, and (2) LVLMs' awareness of this limitation results in higher alignment.

To validate our capacity hypothesis of instruction following ability, we conduct a comparative analysis on different scale models and find that:

As LVLMs scale up, they generally exhibit stronger instruction following capabilities. As shown in Figure 2. For Qwen2.5-VL and DeepSeek-VL2, the Punish method effectively reduces overconfidence in larger models (Qwen2.5-VL-72B, DeepSeek-V12-16B) but shows limited impact on smaller ones ( < 32B Qwen2.5-VL, DeepSeek-VL2-3B). For LLaVA-1.5, the 13B model follows Explain instruction which 7B model not follows, thus Explain improves accuracy in the 13B model.

These phenomena supports our hypothesis: the parameter constraints of small scale LVLMs create a dilemma between visual processing and linguistic comprehension, resulting in degraded language understanding and consequently weaker instruction following ability. In contrast, larger LVLMs allocate more parameters to language processing, maintaining strong language ability while handling multimodal inputs, thus demonstrating stronger instruction following ability.

To validate our accuracy and alignment hypothe-

			Dyn-VQA											
Model	Task	Unc-R.	Acc	Align.	Conser.	Overco.								
Qwen2.5-VL	"V"QA	0.461	0.223	0.578	0.053	<b>0.369</b>								
	VQA	<b>0.782</b>	0.185	<b>0.762</b>	0.102	0.135								
	QA	0.766	<b>0.252</b>	0.700	<b>0.159</b>	0.141								
DeepSeek-VL2	"V"QA	0.227	0.208	0.435	0.000	<b>0.565</b>								
	VQA	<b>0.788</b>	0.146	<b>0.653</b>	<b>0.141</b>	0.207								
	QA	0.545	<b>0.256</b>	0.559	0.121	0.320								

Table 5: The performance of LVLMs under different query modalities, we add text question at the bottom of the image to generate pure image "V"QA query.

sis, we conduct comparative analysis on text-only QA, vision-text VQA, and vision-only "V"QA modality of queries on LVLMs, our results reveal that: 611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

LVLMs exhibit lower accuracy but higher alignment when responding to multimodal VQA queries. As shown in Table 5, both models demonstrate lower accuracy when answering VQA queries that demand cross-modal understanding ability compared to pure text QA and pure image "V"QA queries. Concurrently, they demonstrate increased uncertain-rate and improved confidence performance alignment for these multimodal queries.

These observations support our hypothesis:

1. Limited cross-modal ability: LVLMs struggle to effectively synthesize information across modalities, leading to reduced answering accuracy on multimodal queries compared to unimodal queries.

2. Capability awareness: When encountering challenging multimodal queries, LVLMs exhibit self-awareness of their limited ability through generating more uncertainty responses. This decreases overconfidence and thus improve alignment.

### 6 Conclusion

In this paper, we present a systematic investigation of knowledge boundary perception in LVLMs, assessing this ability through alignment. First, we evaluate three kinds of confidence, and observe that answer consistency-based confidence reaches the highest alignment, whereas verbalized confidence induces overconfidence. We also evaluate several confidence calibration methods, with our results revealing that reasoning elicitation methods improve accuracy and alignment, while our proposed methods show effectiveness. Second, we compare LVLMs with LLMs, and reveal that while LVLMs exhibit lower QA accuracy, they achieve higher alignment, which is attributable to LVLMs' awareness of their multimodal integration ability limitation. We also observe that LVLMs have weaker instruction following ability than LLMs.

752

753

# 651 Limitations

652First, due to dataset constraints, we only compared653LVLMs and LLMs on Dyn-VQA; broader bench-654marks are needed for future validation. Second,655our analysis did not examine internal model states,656leaving internal mechanistic differences in knowl-657edge boundary perception underexplored. Third,658we focused on binary confidence measures; extend-659ing this to continuous confidence scales could yield660finer-grained insights. These limitations highlight661directions for future work on LVLM evaluation and662interpretability.

#### Ethics Statement

663

670

671

681

682

683

687

695

699

In this paper, all the datasets we use are opensource, and the models we employ are either opensource or widely used. Furthermore, the methods we propose do not induce the model to output any harmful information.

#### References

- Ayush Agrawal, Mirac Suzgun, Lester Mackey, and Adam Tauman Kalai. 2024. Do language models know when they're hallucinating references?
- Amos Azaria and Tom Mitchell. 2023. The internal state of an llm knows when it's lying.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. 2025a. Qwen2.5-vl technical report.
- Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou. 2025b. Hallucination of multimodal large language models: A survey.
- Zhuo Chen, Xinyu Wang, Yong Jiang, Zhen Zhang, Xinyu Geng, Pengjun Xie, Fei Huang, and Kewei Tu. 2025. Detecting knowledge boundary of vision large language models by sampling-based inference.
- Nathan Cooper and Torsten Scholak. 2024. Perplexed: Understanding when large language models are confused.
- Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zichong Yang, Kuei-Da Liao, Tianren Gao, Erlong Li, Kun Tang, Zhipeng Cao, Tong Zhou, Ao Liu, Xinrui Yan, Shuqi Mei, Jianguo Cao, Ziran Wang, and Chao Zheng. 2024. A survey on multimodal large language models for autonomous driving. In <u>Proceedings of</u>

the IEEE/CVF Winter Conference on Applications	
of Computer Vision (WACV) Workshops, pages 958-	
979.	

- Shrey Desai and Greg Durrett. 2020. Calibration of pre-trained transformers.
- Tobias Groot and Matias Valdenegro-Toro. 2024. Overconfidence is key: Verbalized uncertainty evaluation in large language and vision-language models.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. 2017a. On calibration of modern neural networks. In <u>International conference on machine</u> learning, pages 1321–1330. PMLR.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. 2017b. On calibration of modern neural networks.
- Yutao Hu, Tianbin Li, Quanfeng Lu, Wenqi Shao, Junjun He, Yu Qiao, and Ping Luo. 2024. Omnimedvqa: A new large-scale comprehensive evaluation benchmark for medical lvlm. In <u>Proceedings of</u> the IEEE/CVF Conference on Computer Vision and <u>Pattern Recognition (CVPR)</u>, pages 22170–22183.
- Bo Jiang, Shaoyu Chen, Bencheng Liao, Xingyu Zhang, Wei Yin, Qian Zhang, Chang Huang, Wenyu Liu, and Xinggang Wang. 2024. Senna: Bridging large vision-language models and end-to-end autonomous driving.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2023. Large language models are zero-shot reasoners.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation.
- Wei Lan, Wenyi Chen, Qingfeng Chen, Shirui Pan, Huiyu Zhou, and Yi Pan. 2024. A survey of hallucination in large visual language models.
- Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian Lu, Chunyan Miao, and Lidong Bing. 2024. Mitigating object hallucinations in large visionlanguage models through visual contrastive decoding. In <u>Proceedings of the IEEE/CVF Conference on</u> <u>Computer Vision and Pattern Recognition (CVPR)</u>, pages 13872–13882.
- Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023. Llavamed: Training a large language-and-vision assistant for biomedicine in one day. In <u>Advances in Neural</u> <u>Information Processing Systems</u>, volume 36, pages 28541–28564. Curran Associates, Inc.
- Jiaming Li, Jiacheng Zhang, Zequn Jie, Lin Ma, and Guanbin Li. 2025a. Mitigating hallucination for large vision language model by inter-modality correlation calibration decoding.

Yangning Li, Yinghui Li, Xinyu Wang, Yong Jiang, Zhen Zhang, Xinran Zheng, Hui Wang, Hai-Tao Zheng, Fei Huang, Jingren Zhou, and Philip S. Yu. 2025b. Benchmarking multimodal retrieval augmented generation with dynamic vqa dataset and self-adaptive planning agent.

754

755

763

765

773

777

778

779

781

782

784

787

788

790

795

796

799

804

805

- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Teaching models to express their uncertainty in words. arXiv preprint arXiv:2205.14334.
- Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiutian Zhao, Ke Wang, Liping Hou, Rongjun Li, and Wei Peng. 2024a. A survey on hallucination in large vision-language models.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024b. Improved baselines with visual instruction tuning.
- Huan Ma, Jingdong Chen, Guangyu Wang, and Changqing Zhang. 2025. Estimating llm uncertainty with logits.
- Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. 2023a. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models.
- Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023b. Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models. arXiv preprint arXiv:2303.08896.
- Viktor Moskvoretskii, Maria Lysyuk, Mikhail Salnikov, Nikolay Ivanov, Sergey Pletenev, Daria Galimzianova, Nikita Krayko, Vasily Konovalov, Irina Nikishina, and Alexander Panchenko. 2025. Adaptive retrieval without self-knowledge? bringing uncertainty back home.
- Shiyu Ni, Keping Bi, Jiafeng Guo, and Xueqi Cheng.
   2024. When do llms need retrieval augmentation? mitigating llms' overconfidence helps retrieval augmentation.
- Shiyu Ni, Keping Bi, Jiafeng Guo, Lulu Yu, Baolong Bi, and Xueqi Cheng. 2025. Towards fully exploiting llm internal states to enhance knowledge boundary perception.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and et al. Red Avila. 2024. Gpt-4 technical report.
- Pranab Sahoo, Prabhash Meharia, Akash Ghosh, Sriparna Saha, Vinija Jain, and Aman Chadha. 2024. A comprehensive survey of hallucination in large language, image, video and audio foundation models.
- Fengfei Sun, Ningke Li, Kailong Wang, and Lorenz Goette. 2025. Large language models are overconfident and amplify human bias.

Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D. Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback.

807

808

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

847

848

849

850

851

852

853

854

855

856

857

858

- Xintong Wang, Jingheng Pan, Liang Ding, and Chris Biemann. 2024a. Mitigating hallucinations in large vision-language models with instruction contrastive decoding.
- Yuhao Wang, Zhiyuan Zhu, Heyang Liu, Yusheng Liao, Hongcheng Liu, Yanfeng Wang, and Yu Wang. 2024b. Drawing the line: Enhancing trustworthiness of mllms through the power of refusal.
- Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, Zhenda Xie, Yu Wu, Kai Hu, Jiawei Wang, Yaofeng Sun, Yukun Li, Yishi Piao, Kang Guan, Aixin Liu, Xin Xie, Yuxiang You, Kai Dong, Xingkai Yu, Haowei Zhang, Liang Zhao, Yisong Wang, and Chong Ruan. 2024. Deepseekvl2: Mixture-of-experts vision-language models for advanced multimodal understanding.
- Wenyi Xiao, Ziwei Huang, Leilei Gan, Wanggui He, Haoyuan Li, Zhelun Yu, Fangxun Shu, Hao Jiang, and Linchao Zhu. 2025. Detecting and mitigating hallucination in large vision language models via fine-grained ai feedback. <u>Proceedings of the AAAI</u> <u>Conference on Artificial Intelligence</u>, 39(24):25543– 25551.
- Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2024. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms.
- Adam Yang, Chen Chen, and Konstantinos Pitas. 2024a. Just rephrase it! uncertainty estimation in closedsource language models via multiple rephrased queries.
- Yuqing Yang, Ethan Chern, Xipeng Qiu, Graham Neubig, and Pengfei Liu. 2024b. Alignment for honesty.
- Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. 2023. Do large language models know what they don't know?
- Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, Yu Su, Wenhu Chen, and Graham Neubig. 2024. Mmmu-pro: A more robust multidiscipline multimodal understanding benchmark.
- Dylan Zhang, Xuchao Zhang, Chetan Bansal, Pedro Las-Casas, Rodrigo Fonseca, and Saravan Rajmohan. 2023. Pace-lm: Prompting and augmentation for calibrated confidence estimation with gpt-4 in cloud incident root cause analysis.

A.1 Prompts       94         Nyang Zhou, Chenhang Cui, Jachong Yoo, Linjung Huaku Yoo, 2024, Analyzing and mitigating an	Jiaxin Zhang, Zhuohang Li, Kamalika Das, Bradley A. Malin, and Sricharan Kumar 2024, Sao3: Paliable	A Appendix	873
A.1.1 Single Step Verbalization Based       69         Yiyang Zhan Deng, Chelsea Finn, Mohit Bansal, and Huaxii Yao. 2024. Analyzing and mitigating other hallhacination in large vision-language models.       Yamilla. Answer the question based on your inter- nal knowledge and the image. If you are sure the auswer is accurate and correct, please say "certain".         Yuke Zhu, Oliver Groth. Michael Bernstein, and Li Feit-Fourte. Conquier Vision and Pattern Recognition (CVPR).       Yanita. Answer the question based on your internal knowledge and the image. Analyse step by step. If you are sure the answer. If you are sure the asswer is accurate and correct, please say "certain".         Question: [Question] Answer:       Image. Coll. Answer the question based on your internal knowledge and the image. Analyse step by step. If you are sure the answer. If you are sure not confident with the answer. If you are sure that involvedge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "certain".         Question: [Question] Answer:       Ting-Coll. Answer the question based on your inter- nal knowledge and the image. First, describe the image, then answer is accurate and correct, please say "cer- tain" after the answer. If you are sure the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the answer is accurate and correct, please say "certain" Question: [Question] Answer:         Fundsh. Answer the question based on your inter- nal knowledge and the image, explain why you give this answer. If you are not confident with the answere; please say "certain" Question: [Question] Answer:<	hallucination detection in black-box language models	A.1 Prompts	874
<ul> <li>Tryang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Thung, Xhung, Zhun Dang, Chesa Finn, Mohi Buai, and Huaxiu Yao. 2024. Analyzing and mitigating object hallucination in large vision-inguage models of the answer the question based on your internal theorem of the constraints. The computer vision and Pattern Recognition (CVPR).</li> <li>Yoke Zhu, Oliver Groth, Michael Bernstein, and Li Feirél. 200 (Constraints). The computer vision and Pattern Recognition (CVPR).</li> <li>Coff. Answer the question based on your internal knowledge and the image. <i>Life you are not confident with the answer</i>. Journal of the computer vision and Pattern Recognition (CVPR).</li> <li>Coff. Answer the question based on your internal knowledge and the image. <i>Analyse step by step</i>. Jf you are sure the answer is accurate and correct, please say "uncertain". Question: [Question] Answer:</li> <li>Coff. Answer the question based on your internal knowledge and the image. <i>First, describe the image, then analyse step by step</i>. Jf you are sure the answer is accurate and correct, please say "uncertain". Question: [Question] Answer:</li> <li>Ing-Coff. Answer the question based on your internal knowledge and the image. <i>First, describe the image, then analyse step by step</i>. Jf you are sure the answer is accurate and correct, please say "certain" after the answer. Jf you are sure the answer is accurate and correct, please say "certain". Question: [Question] Answer:</li> <li>Punish. Answer the question based on your internal knowledge and the image. <i>First, describe the image, then analyse step by step</i>. Jf you are sure the answer is accurate and correct, please say "certain" after the answer. Jf you are sure the answer is accurate and correct, please say "certain" after the answer. Jf you are sure the answer is accurate and correct, please say "certain" after the answer. Jf you are sure the answer is accurate and correct, please say "certain" after the answer. Jf you are sure the answer is accurate and correct, please say "certain" after the a</li></ul>	via semanic-aware cross-check consistency.	A.1.1 Single Step Verbalization Based	875
Yang, Jain Deng, Outersal Pring, and mitigating objet       Yanilla. Answer the question based on your inter- answer is accurate and correct, please say "certain"       979         Yuke Zho, Oliver Gorb, Michael Bernstein, and to massering in images. In Proceedings of the Recegation (CVPR).       You are not confident with the answer is accurate and correct, please say "certain".       983         Col.       Answer:       983         Question: [Question]       985       99       959         Answer:       983       986       993       987         Question: [Question]       986       993       987       993       987         Question: [Question]       986       993       987       9933       993       993 <th>Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun</th> <th>Prompts</th> <th>876</th>	Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun	Prompts	876
hallucination in large vision-language models.Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-F. 2016. Vision and Patter Recognition (CVPR).BifER, Conference on Computer Vision and Patter Recognition (CVPR).Coll. Answer The answer, Jeause say "uncertain". Question: [Question] Answer:Coll. Answer the question based on your internal knowledge and the image. Anguse step by step. If you are surve the answer. Joyan are not confident with the answer. Joyan are not confident with the answer is accurate and correct, please say "certain". Question: [Question] Answer:Ing-CoT. Answer the question based on your in- ternal knowledge and the image. First, describe the image, then analyse step by step. J step. Answer:Ing-CoT. Answer the question based on your in- ternal knowledge and the image. First, describe the image. Step by step. J step. Answer:Ing-CoT. Answer the question based on your inter- nal knowledge and the image. First, describe the image. Step by step. J step. Step are not confident with the answer: Joyan are not confident with the answer: Joyan are not confident with the answer: Joyan are not confident with the answer: Step are step.Ing-CoT. Answer the question based on your inter- nal knowledge and the image. First, describe the image. Step by step.Ing-CoT. Answer the question based on your inter- nal knowledge and the image.	Huaxin Yao 2024 Analyzing and mitigating object	Vanilla. Answer the question based on your inter-	877
Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. 2016. Visual7w: Grounded que- tion answering in images. In Proceedings of the ITFJF: Conference on Computer Vision and Pattern Recognition (CVPR).       answer is accurate and correct, please say "uncertain". Question: [Question]       881         Colf. Answer the question based on your internal knowledge and the image. Analyse step by step. If you are sure the answer is accurate and correct, please say "certain" (after the answer. If you are not confident with the answer: Please say "certain" (after the answer. If you are not conference on computer Vision and Pattern Recognition (CVPR).       881         Colf. Answer the question based on your internal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain". (Question] Answer:       882         Img-Colf. Answer the question based on your inter- nal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the answer. Please say "uncertain". Question: [Question] Answer:       883         Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain" offer the answer. If you are sure the answer: is not right but you say "certain". Question: [Question] Answer:       893         Punish. Answer the question based on your inter- nal knowledge and the image. Analyse set on your inter- nal knowledge and the image. Tyou are sure the answer. If you are not confident with the answer. Please say "uncertain". Question: [Question] Answer:       893	hallucination in large vision-language models.	nal knowledge and the image. If you are sure the	878
The Price And, Oriver Ground, The Ground and Particle And Partend Parte And Particle And Particle And Particle And Pa	Value 7km Oliver Creth Michael Demotrie and	answer is accurate and correct, please say "certain"	879
<ul> <li>answer, please say "uncertain".</li> <li>Question: [Question]</li> <li>Answer:</li> <li>Col. Answer the question based on your internal knowledge and the image. Analyse step by step. If you are sure the answer is accurate and correct, please say "certain" after the answer. If you are not confident with the answer, accurate and correct, please say "certain".</li> <li>Question: [Question]</li> <li>Answer:</li> <li>Image, then analyse is accurate and correct, please say "certain".</li> <li>Question: [Question]</li> <li>Answer:</li> <li>Prinsh. Answer the question based on your internal knowledge and the image. First, describe the image, then analyse is accurate and correct, please say "certain".</li> <li>Question: [Question]</li> <li>Answer:</li> <li>Thing-Col. Answer the question based on your internal file answer is accurate and correct, please say "certain" after the answer is accurate and correct, please say "certain" after the answer. If you are sure the answer is accurate and correct, please say "certain".</li> <li>Question: [Question]</li> <li>Answer:</li> <li>Punish. Answer the question based on your internal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain" after the answer. If you are not confident with the answer is accurate and correct, please say "certain".</li> <li>Question: [Question]</li> <li>Answer:</li> <li>Explain. Answer the question based on your internal knowledge and the image. Erg you are sure the answer is not right but you say "certain".</li> <li>Question: [Question]</li> <li>Answer:</li> <li>Question: [Question]</li> <li>Answer:</li> <li>Conditionation and the answer is accurate and correct, please say "certain" after the answer. If you are sure the answer is not right but you say "certain".</li> <li>Question: [Question]</li> <li>Answer:</li> <li>Conditionation anonon confident with the answer.</li> <li>Courcertain".</li></ul>	Li Fei-Fei 2016 Visual <sup>7</sup> w: Grounded ques-	after the answer. If you are not confident with the	880
IEEE Conference on Computer Vision and Pattern Recognition (CVPR).Question: [Question] Answer:483CoT. Answer the question based on your internal knowledge and the image. Analyse step by step. If you are sure the answer is accurate and correct, please say "certain" after the answer. If you are not confident with the answer, please say "uncer- tain". Question: [Question] Answer:883Img-CoT. Answer the question based on your in- ternal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "uncer- tain". Question: [Question] Answer:883Img-CoT. Answer the question based on your in- ternal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the answer:Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain" difer the answer. If you are not confident with the answer is not right but you say "certain". Question: [Question] Answer:Answer:707Explain. Answer the question based on your inter- nal knowledge and the image. explain why you give this answer. If you are sure the answer is accurate and correct, please say "certain". Sourcentain". Question: [Question] Answer:Masseer:707Explain. Answer the question based on your inter- nal knowledge and the image. explain why you give this answer. If you are not confident with the answer, if you	tion answering in images. In Proceedings of the	answer, please say "uncertain".	881
Recognition (CVPR).Answer:683CoT. Answer the question based on your internal knowledge and the image. Analyse step by step. If you are sure the answer is accurate and correct, please say "certain" after the answer. If you are not confident with the answer, please say "uncer- tain".683Question: [Question] Answer:683Question: [Question] Answer:683Prime Cott. Answer the question based on your in- ternal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain" after the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the answer: gease say "uncertain". Question: [Question] Answer:683Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer:693Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer: Question: [Question] Answer:693Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer: please say "uncertain". Question: [Question] Answer:693Punish. Answer the question based on your inter- nal knowledge and the image. explain why you give this answer. [you are sure the answer is accurate and correct, please say "certain" offer the answer is not right but you say "certain". Question: [Question] Answer:693Question: [Question] Answer:693693Question: [Question] Answer:693Question: [Question] Answer:693Question: [Question] Answer:693 </th <th>IEEE Conference on Computer Vision and Pattern</th> <th>Question: [Question]</th> <th>882</th>	IEEE Conference on Computer Vision and Pattern	Question: [Question]	882
CoT. Answer the question based on your internal knowledge and the image. Analyse step by step. If you are sure the answer is accurate and correct, please say "certain" after the answer. If you are not confident with the answer, please say "uncer- tain". Question: [Question] Answer:Img-CoT. Answer the question based on your in- ternal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain" (der the answer, fly ou are not confident with the answer is accurate and correct, please say "cer- tain" (der the answer the question based on your inter- tain" (der the answer the question based on your inter- tains). Answer:Punish. Answer the question based on your inter- tains. Answer the question based on your inter- tains. If you are not confident with the answer:Busish. Answer the question based on your inter- after the answer. If you are not confident with the answer:Cuestion: [Question] Answer:Question: [Question] Answer:Question: [Question] Answer:Question: [Question] Answer:Question: [Question] Answer:Question: [Question] Answer:Question: [Question] Answer:Question: [Question] Answer:Question: [Question] Answer:Question: [Question] Answer:PromptsCother answer. By ou are sure the answer. By ou are sure the answer. So on triphent with the answer:Question: [Question] Answer:Question: [Question] Answer:Question: [Question] Answer:Question: [Question] Answer:Question: [Question] Answer:Questio	Recognition (CVPR).	Answer:	883
knowledge and the image. Analyse step by step. If you are sure the answer is accurate and correct, please say "certain" difer the answer. If you are not confident with the answer, please say "uncer- tain". Question: [Question] Answer:889Ing-CoT. Answer the question based on your in- ternal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain". Question: [Question] Answer:Punish. Answer the question based on your inter- ral "after the answer. If you are not confident with the answer, please say "uncertain". Question: [Question] Answer:Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer:Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer:Punish. Answer the question based on your inter- nal knowledge and the image. Toy our as sure the answer: is accurate and correct, please say "ecretain" apter the answer. If you are not confident with the answer: so accurate and correct, please say "uncertain". Question: [Question] Answer:Bis answer, if you are not confident with the answer, please say "uncertain". You will be punish- ed if the answer. If you are not confident with the answer; we please say "uncertain". You are not confident with the answer;		<b>CoT.</b> Answer the question based on your internal	884
you are sure the answer is accurate and correct, please say "certain" (fler the answer. If you are not confident with the answer, please say "uncer- tain". Question: [Question] Answer: Ing-CoT. Answer the question based on your in- ternal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the answer: Question: [Question] Answer: Punish. Answer the question based on your inter- nal knowledge and the image. Jyou are sure the answer is accurate and correct, please say "certain" after the answer. If you are not confident with the answer is accurate and correct, please say "certain" after the answer. If you are not confident with the answer is accurate and correct, please say "certain" after the answer. If you are not confident with the answer is accurate and correct, please say "certain" after the answer is not right but you say "certain". Question: [Question] Answer: Question: [Question] Answer:		knowledge and the image. Analyse step by step. If	885
please say "certain" after the answer. If you are not confident with the answer, please say "uncer- tain".887 not confident with the answer, please say "uncer- tain".Question: [Question] Answer:889Question: [Question] Answer:891Inge-CoT. Answer the question based on your in- ternal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the answer please say "uncertain". Question: [Question] Answer:Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain" guestion: [Question] Answer:Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain" guestion: [Question] Answer:Punish. Answer the question based on your inter- nal knowledge and the image, explain why you give this answer:907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give this answer:909Answer:909Answer:909Answer:909Answer:901902903904905905906907908909909909909909909909909909909909909 </th <th></th> <th>you are sure the answer is accurate and correct,</th> <th>886</th>		you are sure the answer is accurate and correct,	886
not confident with the answer, please say "uncertain".88Question: [Question]99Answer:91Ing-CoT. Answer the question based on your internal knowledge and the image. First, describe the93image, then analyse step by step. If you are sure94the answer is accurate and correct, please say "cer-95tain" affer the answer. If you are not confident with96the answer is accurate and correct, please say "certain"97Question: [Question]98Answer:99Punish. Answer the question based on your internal knowledge and the image. If you are sure the90answer is accurate and correct, please say "certain"90after the answer is accurate and correct, please say "certain"90after the answer is not right but you say "certain".90Question: [Question]90Answer:90answer:90Question: [Question]90Answer:90Question: [Question]90Answer:90Question: [Question]90Answer:90Question: [Question]90Answer:90Question: [Question]90Answer:90Answer:91Question: [Question]90Answer:91Question: [Question]91Answer:91Answer:91Answer:91Answer:92You are sure the answer, please say "certain"9392 <tr< th=""><th></th><th>please say "certain" after the answer. If you are</th><th>887</th></tr<>		please say "certain" after the answer. If you are	887
tain".899Question: [Question]891Answer:891Ing-CoT. Answer the question based on your in- ternal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the enswer:Question: [Question]893Answer:893Punish. Answer the question based on your inter- nal knowledge and the image, If you are sure the 		not confident with the answer, please say "uncer-	888
Question: [Question] Answer:890Img-CoT. Answer the question based on your in- ternal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain" dfer the answer. If you are not confident with the answer; please say "uncertain". Question: [Question] Answer:Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain" qfier the answer is you are not confident with the answer:Question: [Question] Answer:Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer: please say "uncertain". You will be punish- ed if the answer is not right but you say "certain". Question: [Question] Answer:B07Question: [Question] Answer:B07Question: [Question] Answer:B07Question: [Question] Answer:B07D08Answer:B07B07Question: [Question] Answer:B07B08Answer:B07B08B08B09<		tain".	889
Answer:891Ing-CoT. Answer the question based on your internal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the answer. If you are not confident with the answer. If you are not confident with the answer.Question: [Question]893Answer:899Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain" after the answer. Jou are not confident with the answer is accurate and correct, please say "certain" ou are not confident with the answer is accurate and correct, please say "certain" ou say "certain". Question: [Question] Answer:901Answer:902901Answer:903Question: [Question] Answer:905Question: [Question] Answer:905Question: [Question] Answer:905Question: [Question] Answer:905Question: [Question] Answer:905Question: [Question] Answer:905Question: [Question] Answer:903Answer:903904 Uncertain".903905903904 Uncertain".904905904906 Answer:905907905908904909904909905909905909906909907909908901909901<		Question: [Question]	890
Img-CoT. Answer the question based on your internal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- atin" after the answer. If you are not confident with the answer.Question: [Question]899Answer:899Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain" of the answer. If you are not confident with the answer is accurate and correct, please say "certain" of the answer is accurate and correct, please say "certain" of the answer is not confident with the answer is not confident with the answer.901 901 901 905 905 906 906 907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give this answer. If you are not confident with the answer:903 905 905 906 905 906 905 906 905 906 905 907906 906 905 906 907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give this answer. If you are not confident with the answer: 907 907 901 901 901 901 901 901 901 901 901 902 903 903 903 903 903 903 904 907 906 907 907 907 908 909 909 909 904 909 909 909 909 904 <br< th=""><th></th><th>Answer:</th><th>891</th></br<>		Answer:	891
ternal knowledge and the image. First, describe the image, then analyse step by step. If you are sure the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the answer, please say "uncertain". Question: [Question] Answer:997Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain" after the answer is not right but you say "certain" off the answer is not right but you say "certain". Question: [Question] Answer:900Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain" off the answer is not right but you say "certain" off the answer is not right but you say "certain". 903 Question: [Question] Answer:907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give please say "certain". 907906Answer:907907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give please say "certain" after the answer: 907907Explain. Answer is not right but you say "certain". 905 Question: [Question] Answer:909909901 His answer, If you are sure the answer is accurate and correct, please say "certain" after the answer. 911 If you are not confident with the answer, securate and correct, please say "certain" after the answer. 913914 Answer:915915A.1.2 Double Step Verbalization Based Prompts917916 917 917917917914 400918 <th></th> <th><b>Img-CoT.</b> Answer the question based on your in-</th> <th>892</th>		<b>Img-CoT.</b> Answer the question based on your in-	892
image, then analyse step by step.If you are surethe enswer is accurate and correct, please say "cer-tain " after the answer. If you are not confident withthe answer, please say "uncertain".Question:Question:Question:Runish.Answer:Base say "certain"Answer:Base say "certain"Base say "uncertain".Base say "certain"Base say "uncertain".Base		ternal knowledge and the image. First, describe the	893
the answer is accurate and correct, please say "cer- tain" after the answer. If you are not confident with the answer, please say "uncertain".897 Question: [Question]Question: [Question]Answer:Punish. Answer the question based on your inter- 		image, then analyse step by step. If you are sure	894
tain" after the answer. If you are not confident with the answer, please say "uncertain".997 Question: [Question]993 Answer:Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain"900 answer is accurate and correct, please say "certain" 903 after the answer. If you are sure the answer is not right but you say "certain". 904 ed if the answer is not right but you say "certain". 905 Question: [Question] Answer:906 907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give please say "uncertain". You will be punish- 904 ed if the answer is not right but you say "certain". 905 Question: [Question] Answer:907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give pli you are not confident with the answer. 907903 904 904 905 906 906 906 907Answer:907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give 909 900 900 900 900 900 900 90		the answer is accurate and correct, please say "cer-	895
the answer, please say "uncertain".697 Question: [Question]697 Question: [Question]697 Question: [Question]Answer:899Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain"900 after the answer. If you are not confident with the 903 answer, please say "uncertain". You will be punish- 904 ed if the answer is not right but you say "certain". 905 Question: [Question]906 906 906 Answer:Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give this answer. If you are sure the answer is accurate and correct, please say "certain" after the answer. 911 H you are not confident with the answer. 911 H you are not confident with the answer. 911 Answer:903Cuestion: [Question]904 905 907905 907Answer:907907Answer:907Duble Step Verbalization Based Prompts916 917For the double step verbalization based methods, we first prompt the model to generate answer, then919		tain" after the answer. If you are not confident with	896
Question: [Question]899Answer:899Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain"900after the answer. If you are not confident with the answer, please say "uncertain". You will be punish- ed if the answer is not right but you say "certain". 905 Question: [Question]906Answer:907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give this answer. If you are sure the answer is accurate and correct, please say "certain" after the answer. 911 If you are not confident with the answer. 911 Question: [Question]905Answer:913 Question: f you are sure the answer, please say "uncertain".914 Answer:Answer:915915A.1.2 Double Step Verbalization Based Prompts916 917For the double step verbalization based methods, we first prompt the model to generate answer, then919		the answer, please say "uncertain".	897
Answer:899Punish. Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain" after the answer. If you are not confident with the answer, please say "uncertain". You will be punish- ed if the answer is not right but you say "certain". 906 Answer:907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give this answer. If you are not confident with the and correct, please say "certain". 907906 906 906 907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give this answer. If you are not confident with the answer. If you are not confident with the answer. If you are not confident with the answer. 909 111 If you are not confident with the answer. 911 If you are not confident with the answer. 913 Question: [Question] Answer:916 917 917 916For the double Step Verbalization based methods, we first prompt the model to generate answer, then919 917		Question: [Question]	898
Punish.Answer the question based on your inter- nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain"902 after the answer. If you are not confident with the 903 answer, please say "uncertain".904 904 ed if the answer is not right but you say "certain". 905 Question: [Question] Answer:906 906 907Explain.Answer the question based on your inter- nal knowledge and the image, explain why you give this answer. If you are sure the answer is accurate and correct, please say "certain" after the answer. 907906 907Explain.Answer the question based on your inter- nal knowledge and the image, explain why you give this answer. If you are sure the answer is accurate and correct, please say "certain" after the answer. 911 If you are not confident with the answer, please say "uncertain". 913 Question: [Question] Answer:916 917For the double Step Verbalization Based Prompts916 917917For the double step verbalization based methods, we first prompt the model to generate answer, then918		Answer:	899
nal knowledge and the image. If you are sure the answer is accurate and correct, please say "certain"902 after the answer. If you are not confident with the 903 answer, please say "uncertain". You will be punish- 904 ed if the answer is not right but you say "certain". 905 Question: [Question]906 906 906 907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give and correct, please say "certain" and correct, please say "certain" and correct, please say "certain" and correct, please say "certain" after the answer. 911 If you are not confident with the answer. 911 If you are not confident with the answer. 913 Question: [Question]913 913 913 914 915A.1.2 Double Step Verbalization Based Prompts916 917 917916For the double step verbalization based methods, we first prompt the model to generate answer, then919		<b>Punish.</b> Answer the question based on your inter-	900
answer is accurate and correct, please say "certain"902after the answer. If you are not confident with the903answer, please say "uncertain". You will be punish-904ed if the answer is not right but you say "certain".905Question: [Question]906Answer:907Explain. Answer the question based on your inter-908nal knowledge and the image, explain why you give909this answer.907Explain. Answer the question based on your inter-908nal knowledge and the image, explain why you give909this answer.910and correct, please say "certain" after the answer.911If you are not confident with the answer, please say912"uncertain".913Question: [Question]914Answer:915A.1.2 Double Step Verbalization Based916Prompts917For the double step verbalization based methods,918we first prompt the model to generate answer, then919		nal knowledge and the image. If you are sure the	901
after the answer. If you are not confident with the answer, please say "uncertain". You will be punish- ed if the answer is not right but you say "certain".904 904 904 905 906 906 906 906 Answer:906 906 906 907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give 909 this answer. If you are sure the answer is accurate 910 and correct, please say "certain" after the answer.906 909 909 911 911 919 912 912 912 911 913 Question: [Question] 914 Answer:908 918 919A.1.2 Pouble Step Verbalization based methods, we first prompt the model to generate answer, then918 919		answer is accurate and correct, please say "certain"	902
answer, please say "uncertain". You will be punish- ed if the answer is not right but you say "certain".904ed if the answer is not right but you say "certain".905Question: [Question]906Answer:907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give pose909this answer.1f you are sure the answer is accurate and correct, please say "certain" after the answer.911If you are not confident with the answer, please say912 "uncertain".913Question: [Question]914Answer:915A.1.2 Double Step Verbalization Based Prompts916For the double step verbalization based methods, we first prompt the model to generate answer, then919		after the answer. If you are not confident with the	903
ed if the answer is not right but you say "certain".905Question: [Question]906Answer:907Explain. Answer the question based on your internal knowledge and the image, explain why you give909this answer.If you are sure the answer is accurateand correct, please say "certain" after the answer.911If you are not confident with the answer, please say912"uncertain".913Question: [Question]914Answer:915A.1.2 Double Step Verbalization Based916Prompts917For the double step verbalization based methods, we first prompt the model to generate answer, then919		answer, please say "uncertain". You will be punish-	904
Question: [Question]906Answer:907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give pose this answer. If you are sure the answer is accurate and correct, please say "certain" after the answer. 911 If you are not confident with the answer, please say "uncertain". 913 Question: [Question] Answer:913A.1.2 PromptsDouble Step Verbalization Based Prompts916 917For the double step verbalization based methods, we first prompt the model to generate answer, then918		ed if the answer is not right but you say "certain".	905
Answer:907Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give nal knowledge and the image, explain why you give this answer. If you are sure the answer is accurate and correct, please say "certain" after the answer.909 900 913 914 914 914 914 914 915 915916 917 917 917 917 917 917 917 917 918 919 919916 917 917 917 918 919		Question: [Question]	906
Explain. Answer the question based on your inter- nal knowledge and the image, explain why you give 909 this answer. If you are sure the answer is accurate and correct, please say "certain" after the answer.909 909 909 910 911 911 912 913 913 Question: [Question] Answer:908 909 909 910 913 914 915A.1.2 Double Step Verbalization Based Prompts916 917 917For the double step verbalization based methods, we first prompt the model to generate answer, then918 919		Answer:	907
nal knowledge and the image, explain why you give this answer. If you are sure the answer is accurate and correct, please say "certain" after the answer.909this answer.If you are sure the answer is accurate and correct, please say "certain" after the answer.910If you are not confident with the answer, please say912"uncertain".913Question: [Question]914Answer:915A.1.2 Double Step Verbalization Based Prompts916917For the double step verbalization based methods, we first prompt the model to generate answer, then919		Explain. Answer the question based on your inter-	908
this answer.If you are sure the answer is accurate910and correct, please say "certain" after the answer.911If you are not confident with the answer, please say912"uncertain".913Question: [Question]914Answer:915A.1.2Double Step Verbalization Based Prompts916917For the double step verbalization based methods, we first prompt the model to generate answer, then919		nal knowledge and the image, <b>explain why you give</b>	909
and correct, please say "certain" after the answer.911If you are not confident with the answer, please say912"uncertain".913Question: [Question]914Answer:915A.1.2 Double Step Verbalization Based Prompts916917For the double step verbalization based methods, we first prompt the model to generate answer, then918		this answer. If you are sure the answer is accurate	910
If you are not confident with the answer, please say912"uncertain".913Question: [Question]914Answer:915A.1.2 Double Step Verbalization Based916Prompts917For the double step verbalization based methods, we first prompt the model to generate answer, then918		and correct, please say "certain" after the answer.	911
"uncertain".913Question: [Question]914Answer:915A.1.2 Double Step Verbalization Based916Prompts917For the double step verbalization based methods,918we first prompt the model to generate answer, then919		If you are not confident with the answer, please say	912
Question: [Question]914Answer:915A.1.2 Double Step Verbalization Based916Prompts917For the double step verbalization based methods,918we first prompt the model to generate answer, then919		"uncertain".	913
Answer:915A.1.2Double Step Verbalization Based Prompts916917917For the double step verbalization based methods, we first prompt the model to generate answer, then918919		Question: [Question]	914
A.1.2Double Step Verbalization Based916Prompts917For the double step verbalization based methods, we first prompt the model to generate answer, then918		Answer:	915
Prompts917For the double step verbalization based methods, we first prompt the model to generate answer, then918		A.1.2 Double Step Verbalization Based	916
For the double step verbalization based methods,918we first prompt the model to generate answer, then919		Prompts	917
we first prompt the model to generate answer, then 919		For the double step verbalization based methods,	918
		we first prompt the model to generate answer, then	919

	900
	957
	958
n	959
	960
S	961
	962
lel	963
1-	964
)-	965
у	966
d	967
1-	968
d.	969
e-	970
1-	971
);	972
<b>1</b> -	973
e	974
У	975
а	976
	977
e	978
У	979
d	980
n	981
h	982
n	983
1-	984
8	985
y n	900
11	907
	900
	989
1-	990
s	991
S	992
	993
	994
s'	995
C-	996
<b>1</b> -	997
	998
or	999
n	1000
1-	1001
vl-	1002
<b>a</b> -	1003
	1004

- prompt the model to give its confidence in the sec-ond round chat.
- 922 First Round Answer Generation. Answer the
  923 question based on your internal knowledge and the
  924 image.
  925 Question: [Question]
  926 Answer:

927 Self-Judging. If you are sure your previous an928 swer is accurate and correct, please say "certain",
929 If you are not confident with the answer, please say
930 "uncertain".

931 CoT. If you are sure your previous answer is accurate and correct, please say "certain", If you
932 are not confident with the answer, please say "uncertain". Analyse step by step, then provide Your
935 judgement.

936 Challenge. *I don't think your answer is right*, if
937 you still think your answer is right, please say "cer938 atin". Otherwise, say "uncertain".

Punish. If you are sure your previous answer is accurate and correct, please say "certain", If you are not confident with the answer, please say "uncertain". You will be punished if the answer is not right but you say "certain".

Probability+Threshold. Provide the probability
that your answer is correct (0.0 to 1.0). Give ONLY
the probability, no other words or explanation.

947 A.1.3 Answer Consistency Based Prompts

948**Rephrasing.** Based on the Following question,949generate [number of semantical equivalent ques-950tions] semantically equivalent questions. your out-951put should be a list of strings and add a sequnce952number with a dot at the start of each output ques-953tion, like [1. "question1", 2. "question2",...].

- 954 Question: [The original question]
- 955 Semantically equivalent questions:

#### A.2 LVLMs' Knowledge Boundary Perception Ability

#### A.2.1 Implementation Details

In this section, we provide a detailed introduction to our implementation details.

For content generation, we mainly utilize APIs to generate answers.

For verbalization based methods, we set the model temperature to 0 and set a fixed seed to obtain highquality and relatively consistent responses. Notably, Probability Threshold method is exclusively employed in a double round form because we find some of the models struggle to generate both continuous probabilities and answers in a single round.

For the consistency based methods, we implemente a two-phase generation protocol: First, generating a reference answer with temperature = 0; Then sampling 10 variant answers with temperature = 1.0, with semantic equivalence between the basic answer and sampled answers evaluated by Qwen2.5-0.5B. With this process, we can get a consistency score between 0 to 10.

Specifically, for question rephrasing method, we leveraged Qwen2.5-7B to produce semantically equivalent question paraphrases. For the noised image method, we progressively added zero-mean Gaussian noise to the images during sampling, with the standard deviation incrementally increased from 0 in steps of 0.05. And for the cross model consistency method, we computed consistency scores using a combination of four responses generated by the primary model and three responses each from two other reference models.

#### A.2.2 Complete Results

Table 6, Table 7 and Table 8 present the comprehensive performance evaluation of all methods across the three benchmark datasets and three LVLMs employed in our study.

#### A.2.3 Observations and Analysis

We proposed our mainly findings about LVLMs' knowledge boundary perception methods in Section 4.4. Here, we discuss more detailed observations about them.

1) The Explain method improves alignment for both Deepseek-VL2 and Qwen2.5 when tested on the Dyn-VQA and Visual7W datasets. This demonstrates its effectiveness in enhancing LVLMs' knowledge boundary perception when processing relatively simple input questions.

1054

1005

1006

2) The single-step Chain of Thought method effectively enhances alignment, whereas its doublestep counterpart often leads to overconfidence and only marginally improves alignment for Qwen2.5-VL.

3) Both single-step and double-step Punish methods demonstrate limited effectiveness in mitigating overconfidence for Qwen2.5-VL and LLaVA-v1.5, as they fail to properly follow Punish Instructions.

4) Challenge method induces very high uncertainrate in both three models, indicating that LVLMs are easily swayed by the output judgements.

5) For Qwen2.5-VL, rephrasing methods improve alignment on the Dyn-VQA dataset (languagefocused), while the noise image method enhances performance on Visual7W (vision-focused). The combination of these two methods boosts alignment on the MMMU Pro dataset, which requires both language and vision comprehension. This reveals an interesting relationship between perturbation modalities and input query types.

# A.3 Comparing Perception between LVLMs and LLMs

While the main body presents a comparative analysis of single-step verbalization based confidence elicitation methods between LLMs and LVLMs, this section provides an extensive evaluation of: (i) double step verbalization based methods, (ii) answer consistency based methods, and (iii) token probability based method. The results can be found in Table 9. The main observations are as follows.

#### A.3.1 Double Step Verbalization Based Methods

For double step verbalization based methods, the difference in performance between LLM and LVLM varies with the method.

1) For the Self-Judging method, Qwen2.5 exhibits higher alignment than Qwen2.5-VL. In contrast, the LLM counterparts of DeepSeek-VL2 and LLaVA tend to respond with "certain" to nearly all answers, resulting in extremely low consistency. This indicates a severe bias toward overconfident responses in these two LLMs.

2) For the Challenge method, LVLMs demonstrate higher uncertain-rates than LLMs, often approaching to near 1.0. This suggests that LVLMs are more likely to trust external judgments and consequently undermine their own decisions.

3) Under the Double-step Punish method, LLMs outperform LVLMs due to their stronger instruction

following ability, achieving higher consistency and lower overconfidence.

1055

1057

1058

1060

1061

1062

1063

1064

1065

1066

1069

1070

1071

1072

1073

1074

#### A.3.2 Answer Consistency Based Methods

For answer consistency based methods, our observations are as follows:

1) Answer consistency based methods demonstrate superior alignment performance in LVLMs compared to LLMs.

2) DeepSeek-MoE exhibits strong consistency in its generated answers, maintaining high answer uniformity even when the outputs are incorrect. This behavior persists across both random sampling and rephrasing methods, leading to sustained overconfidence and suboptimal alignment performance.

3) The rephrasing strategy shows limited effectiveness in improving alignment metrics across all evaluated models, with the notable exception of Qwen2.5-VL. This observation holds true for both LVLMs and LLMs in our results.

#### A.3.3 Token Probability Based Methods

For the token probability based approach, as shown1075in Table 9, our results reveal that LLMs exhibit1076relatively weaker confidence-accuracy alignment1077compared to LVLMs.1078

			Dyn-VQ	A				Visual7W	V				MMMU P	ro	
method	Unc-R.	Acc	Align.	Conser.	Overco.	Unc-R.	Acc	Align.	Conser.	Overco.	Unc-R.	Acc	Align.	Conser.	Overco.
Vanilla	0.7824	0.1846	0.7623	<u>0.1024</u>	0.1353	0.3260	0.4380	0.5840	0.0900	0.3260	0.3000	0.4564	0.4909	0.1327	0.3764
CoT	0.7276	0.2121	0.7824	0.0786	0.1389	0.3840	<u>0.4920</u>	<u>0.6080</u>	0.1340	0.2580	0.2636	<u>0.6436</u>	0.6818	0.1127	0.2055
Img-CoT	0.6545	0.2048	0.7276	0.0658	<u>0.2066</u>	0.2760	0.5020	0.6060	0.0860	0.3080	<u>0.2800</u>	0.6636	0.7182	0.1127	0.1691
Punish	0.7112	0.1682	0.7112	0.0804	0.2084	0.2880	0.4280	0.5520	0.0820	0.3660	0.2691	0.4564	0.5000	0.1127	0.3873
Explain	0.8282	0.1956	0.8117	0.1060	0.0823	<u>0.3640</u>	0.4740	0.6180	<u>0.1100</u>	0.2720	<u>0.2945</u>	0.5309	0.5782	0.1236	0.2982
Self-Judging	0.1426	0.1883	0.3272	0.0018	0.6709	0.1100	0.4760	0.5500	0.0180	0.4320	0.1882	0.5127	0.5609	0.0701	0.3692
СоТ	0.5210	0.1883	0.6435	0.0329	0.3236	0.1500	0.4760	0.5700	0.0280	0.4020	0.2127	0.5127	0.5255	0.1000	0.3745
Challenge	0.9671	0.1883	0.8080	0.1737	0.0183	0.9800	0.4760	0.5280	0.4640	0.0080	0.9873	0.5127	0.4891	0.5055	0.0055
Punish	0.1426	0.1883	0.3272	0.0018	0.6709	0.0780	0.4760	0.5300	0.0120	0.4580	0.0436	0.5127	0.5164	0.0200	0.4636
Prob-Thr	0.4991	0.1883	0.5960	<u>0.0457</u>	<u>0.3583</u>	<u>0.2140</u>	0.4760	0.5820	<u>0.0540</u>	0.3640	<u>0.4764</u>	0.5127	0.5855	<u>0.2018</u>	0.2127
Random	0.4625	0.1883	0.5448	0.0530	0.4022	0.3020	0.4760	0.5700	0.1040	0.3260	0.4309	0.5127	0.5327	0.2055	0.2618
Noised Img	0.7733	0.1883	0.7313	0.1152	0.1536	0.4920	0.4760	0.6000	0.1840	0.2160	0.3873	0.5127	0.5400	0.1800	0.2800
Rephr	0.9543	0.1883	0.8026	0.1700	0.0274	0.4340	0.4760	0.5660	0.1720	0.2620	0.5655	0.5127	0.5364	0.2709	0.1927
Reph+Nois	0.8958	0.1883	0.7733	0.1554	0.0713	<u>0.4940</u>	0.4760	0.5500	0.2100	0.2400	0.4418	0.5127	0.5509	0.2018	0.2473
Cross Model	<u>0.9469</u>	0.1883	0.8208	<u>0.1572</u>	0.0219	0.5720	0.4760	0.6320	<u>0.2080</u>	0.1600	<u>0.5036</u>	0.5127	0.5800	<u>0.2182</u>	0.2018
PPL Thr	0.8885	0.1993	0.7916	0.1481	0.0603	0.8060	0.4760	0.6020	0.3400	0.0580	0.9436	0.4091	0.6073	0.3727	0.0200

Table 6: The performance of different methods on Qwen2.5-VL-7B-Instruct.

			Dyn-VQ	4				Visual7W	/				MMMU P	ro	
method	Unc-R.	Acc	Align.	Conser.	Overco.	Unc-R.	Acc	Align.	Conser.	Overco.	Unc-R.	Acc	Align.	Conser.	Overco.
Vanilla	0.4899	0.0878	0.5338	0.0219	0.4442	0.0260	0.3920	0.4140	0.0020	0.5840	0.0855	0.2018	0.2509	0.0182	0.7309
СоТ	0.5119	0.0841	<u>0.5375</u>	0.0311	0.4314	0.0220	0.3840	0.3940	0.0060	0.6000	0.1418	0.1545	0.2418	0.0273	0.7309
Img-CoT	0.5265	0.0914	0.5484	0.0347	0.4168	0.0180	0.4000	<u>0.4140</u>	0.0020	0.5840	0.1527	0.2527	0.2964	0.0545	0.6491
Punish	0.4497	0.0951	0.4899	0.0274	0.4826	0.0260	0.3960	0.4180	0.0020	0.5800	0.2291	0.2727	0.3745	0.0636	0.5618
Explain	0.4205	0.0841	0.4534	0.0274	0.5192	0.0100	0.3840	0.3900	0.0020	0.6080	0.0727	0.1709	0.2109	0.0164	<u>0.7727</u>
Self-Judging	0.1718	0.1005	0.2468	0.0128	0.7404	0.0020	0.4200	0.4220	0.0000	0.5780	<u>0.0109</u>	0.3218	0.3327	0.0000	0.6673
СоТ	0.0494	0.1005	0.1463	0.0018	0.8519	0.0000	0.4200	0.4200	0.0000	0.5800	0.0000	0.3218	0.3218	0.0000	0.6782
Challenge	1.0000	0.1005	<del>0.8995</del>	0.1005	0.0000	1.0000	0.4200	0.5800	0.4200	0.0000	1.0000	0.3218	0.6782	0.3218	0.0000
Punish	0.0293	0.1005	0.1298	0.0000	0.8702	0.0000	0.4200	0.4200	0.0000	0.5800	0.0000	0.3218	0.3218	0.0000	0.6782
Prob-Thr	0.8464	0.1005	0.7971	0.0750	0.1280	0.6780	0.4200	0.6140	0.2420	0.1440	0.7964	0.3218	0.6091	0.2545	0.1364
Random	0.9872	0.1005	0.8976	0.0951	0.0073	<u>0.6680</u>	0.4200	0.7080	0.1900	0.1020	0.9745	0.3218	0.6709	0.3127	0.0164
Noised Img	0.9963	0.1005	0.8958	0.1005	0.0037	0.5660	0.4200	0.6740	0.1560	0.1700	0.9836	0.3218	0.6655	0.3200	0.0145
Rephr	0.9981	0.1005	0.8976	0.1005	0.0018	0.6560	0.4200	0.6920	<u>0.1920</u>	0.1160	0.9345	0.3218	0.6672	0.2945	0.0382
Reph+Nois	0.9982	0.1005	0.9013	0.0987	0.0000	0.7020	0.4200	0.6780	0.2220	0.1000	0.9655	0.3218	0.6655	0.3109	0.0236
Cross Model	0.9982	0.1005	<u>0.8976</u>	0.1005	0.0018	0.5320	0.4200	0.6520	0.1500	0.1980	0.9727	0.3218	0.6618	<u>0.3164</u>	0.0218
PPL Thr	0.8903	0.1005	<u>0.8519</u>	0.0695	0.0786	0.5860	0.4200	0.7060	0.1460	0.1380	0.9727	0.3218	0.6800	0.3073	0.0127

Table 7: The performance of different methods on LLaVA-v1.5-7B.

			Dyn-VQ/	A				Visual7W	V				MMMU P	ro	
method	Unc-R.	Acc	Align.	Conser.	Overco.	Unc-R.	Acc	Align.	Conser.	Overco.	Unc-R.	Acc	Align.	Conser.	Overco.
Vanilla	<u>0.7879</u>	0.1463	0.6527	0.1408	0.2066	<u>0.4120</u>	0.1840	0.2820	0.1580	0.5600	0.4091	0.2673	0.2727	0.2018	0.5255
CoT	0.6380	0.1700	0.6362	0.0859	<u>0.2779</u>	0.1780	0.4600	0.5540	0.0420	0.4040	0.0873	<u>0.3509</u>	0.3836	0.0273	0.5891
Img-CoT	0.5356	0.2011	0.6344	0.0512	0.3144	0.0640	0.4960	0.5360	0.0120	0.4520	0.1055	0.4582	0.5236	0.0200	0.4564
Punish	0.8483	0.1609	0.7093	0.1499	0.1407	0.4580	0.2680	0.3500	0.1880	0.4620	0.4782	0.3054	0.3145	0.2345	0.4509
Explain	0.7861	0.1682	<u>0.6984</u>	0.1280	0.1737	0.2780	<u>0.4640</u>	0.5700	0.0860	0.3440	0.2073	0.3382	0.3491	0.0982	0.5527
Self-Judging	0.0018	0.1974	0.1993	0.0000	0.8007	0.0020	0.4760	0.4780	0.0000	0.5220	0.0018	0.4255	0.4236	0.0018	0.5745
СоТ	0.0055	0.1974	0.2029	0.0000	0.7971	0.0000	0.4760	0.4760	0.0000	0.5240	0.0073	0.4255	0.4255	0.0036	0.5709
Challenge	0.9945	0.1974	0.8007	0.1956	0.0037	0.9960	0.4760	0.5240	0.4740	0.0020	0.9309	0.4255	0.5709	0.3927	0.0364
Punish	0.3144	0.1974	0.4936	0.0091	0.4973	0.0620	0.4760	0.5300	0.0040	0.4660	0.0273	0.4255	0.4345	0.0091	0.5564
Prob-Thr	<u>0.7239</u>	0.1974	<u>0.6910</u>	<u>0.1152</u>	0.1938	<u>0.7280</u>	0.4760	0.6060	<u>0.2980</u>	0.0960	<u>0.7473</u>	0.4255	<u>0.5218</u>	<u>0.3127</u>	0.1655
Random	0.9963	0.1974	0.8026	0.1956	0.0018	0.5800	0.4740	0.6460	0.2040	0.1500	0.8509	0.4180	0.6000	0.3345	0.0655
Noised Img	0.9927	0.1974	0.8062	0.1920	0.0018	0.5480	0.4740	0.6300	0.1960	0.1740	0.7418	0.4180	0.5818	0.2891	0.1291
Rephr	0.9689	0.1974	<u>0.8080</u>	0.1792	0.0127	0.4480	0.4740	0.6260	0.1480	0.2260	0.7982	0.4180	0.5764	0.3200	0.1036
Reph+Nois	0.9670	0.1974	0.8099	0.1773	0.0128	0.5140	0.4740	0.6120	0.1880	0.2000	0.7945	0.4180	0.5618	0.3255	0.1127
Cross Model	0.9963	0.1974	0.8062	<u>0.1938</u>	0.0000	0.6080	0.4740	0.6740	0.2040	0.1220	0.8327	0.4180	0.5964	<u>0.3273</u>	0.0764
PPL Thr	0.8958	0.1974	0.7934	0.1499	0.0567	0.5500	0.4780	0.6280	0.2000	0.1720	0.9418	0.4436	0.5345	0.4255	0.0400

Table 8: The performance of different methods on DeepSeek-VL2-16B.

		Qwen2.5							LLaVA1.	5		DeepSeek-VL2				
method	Model Type	Unc-R.	Acc	Align.	Conser.	Overco.	Unc-R.	Acc	Align.	Conser.	Overco.	Unc-R.	Acc	Align.	Conser.	Overco.
Self-Judging	LVLM LLM	0.1426 <b>0.2943</b>	0.1883 0.2998	0.3272 <b>0.5649</b>	0.0018 <b>0.0146</b>	<b>0.6709</b> 0.4205	<b>0.0018</b> 0.0000	0.1974 <b>0.2962</b>	0.1993 0.2962	$0.0000 \\ 0.0000$	<b>0.8007</b> 0.7038	<b>0.1718</b> 0.0000	0.1005 0.2139	<b>0.2468</b> 0.2139	<b>0.0128</b> 0.0000	0.7404 <b>0.7861</b>
СоТ	LVLM LLM	<b>0.5210</b> 0.2925	0.1883 <b>0.2998</b>	<b>0.6435</b> 0.5411	<b>0.0329</b> 0.0256	0.3236 <b>0.4333</b>	0.0055 0.2888	0.1974 <b>0.2962</b>	0.2029 0.5192	0.0000 <b>0.0329</b>	<b>0.7971</b> 0.4479	0.0494 <b>0.2761</b>	0.1005 0.2139	0.1463 <b>0.4680</b>	0.0018 <b>0.0110</b>	<b>0.8519</b> 0.5210
Challenge	LVLM LLM	<b>0.9671</b> 0.7148	0.1883 <b>0.2998</b>	<b>0.8080</b> 0.7514	<b>0.1737</b> 0.1316	0.0183 <b>0.1170</b>	<b>0.9945</b> 0.8684	0.1974 <b>0.2962</b>	<b>0.8007</b> 0.6563	0.1956 <b>0.2541</b>	0.0037 <b>0.0896</b>	<b>1.0000</b> 0.9853	0.1005 <b>0.2139</b>	<b>0.8995</b> 0.7898	0.1005 <b>0.2048</b>	0.0000 <b>0.0055</b>
Punish	LVLM LLM	0.1426 <b>0.5448</b>	0.1883 0.2998	0.3272 <b>0.6910</b>	0.0018 <b>0.0768</b>	<b>0.6709</b> 0.2322	<b>0.3144</b> 0.2852	0.1974 <b>0.2962</b>	0.4936 <b>0.5302</b>	0.0091 0.0256	<b>0.4973</b> 0.4442	0.0293 0.1974	0.1005 <b>0.2139</b>	0.1298 <b>0.4113</b>	$0.0000 \\ 0.0000$	<b>0.8702</b> 0.5887
Prob-Thr	LVLM LLM	0.4991 <b>0.5941</b>	0.1883 <b>0.2998</b>	0.5960 <b>0.6709</b>	0.0457 <b>0.1115</b>	<b>0.3583</b> 0.2175	<b>0.7239</b> 0.1773	0.1974 <b>0.2962</b>	<b>0.6910</b> 0.4333	<b>0.1152</b> 0.0201	0.1938 <b>0.5466</b>	0.8464 <b>0.9963</b>	0.1005 <b>0.2139</b>	<b>0.7971</b> 0.7824	0.0750 <b>0.2139</b>	<b>0.1280</b> 0.0037
Random	LVLM LLM	0.4625 <b>0.9287</b>	0.1883 <b>0.2998</b>	0.5448 <b>0.7203</b>	0.0530 <b>0.2541</b>	<b>0.4022</b> 0.0256	<b>0.9963</b> 0.4863	0.1974 <b>0.2962</b>	<b>0.8026</b> 0.5448	<b>0.1956</b> 0.1188	0.0018 <b>0.3364</b>	<b>0.9872</b> 0.8921	0.1005 <b>0.2139</b>	<b>0.8976</b> 0.7806	0.0951 <b>0.1627</b>	0.0073 <b>0.0567</b>
Rephr	LVLM LLM	<b>0.9543</b> 0.9068	0.1883 0.2998	<b>0.8026</b> 0.7203	0.1700 0.2431	0.0274 <b>0.0366</b>	<b>0.9689</b> 0.4991	0.1974 0.2962	<b>0.8080</b> 0.5539	<b>0.1792</b> 0.1207	0.0127 0.3254	<b>0.9981</b> 0.8757	0.1005 0.2139	<b>0.8976</b> 0.7751	0.1005 0.1572	0.0018 <b>0.0676</b>
PPL Thr	LVLM LLM	<b>0.8885</b> 0.8519	0.1993 <b>0.3217</b>	<b>0.7916</b> 0.7313	0.1481 <b>0.2212</b>	<b>0.0603</b> 0.0475	<b>0.8903</b> 0.7587	0.1005 <b>0.2980</b>	<b>0.8519</b> 0.6837	0.0695 <b>0.1865</b>	0.0786 <b>0.1298</b>	<b>0.8958</b> 0.7458	0.1974 <b>0.2121</b>	<b>0.7934</b> 0.7422	<b>0.1499</b> 0.1079	0.0567 <b>0.1499</b>

Table 9: Performance comparison of double step verbaliztion based methods, consistency based methods and answer consistency based methods on the Dyn-VQA dataset: LVLMs vs. LLMs