VLM²-Bench: A Closer Look at How Well VLMs Implicitly Link Explicit Matching Visual Cues

Anonymous ACL submission

Abstract

Visually linking matching cues is a crucial ability in daily life, such as identifying the same person in multiple photos based on their cues, even without knowing who they are. Despite the extensive knowledge that vision-language models (VLMs) possess, it remains largely unexplored whether they are capable of performing this fundamental task. To address this, we introduce VLM²-Bench, a benchmark designed to assess whether VLMs can Visually Link Matching cues, with 9 subtasks and over 3.000 test cases. Comprehensive evaluation across eight open-source VLMs and GPT-4o, along with further analysis of various languageside and vision-side prompting methods, leads to a total of eight key findings. We identify critical challenges in models' ability to link visual cues, highlighting a significant performance gap where even GPT-40 lags 34.80% behind humans. Based on these insights, we advocate for (i) enhancing core visual capabilities to improve adaptability and reduce reliance on prior knowledge, (ii) establishing clearer principles for integrating language-based reasoning in vision-centric tasks to prevent unnecessary biases, and (iii) shifting vision-text training paradigms toward fostering models' ability to independently structure and infer relationships among visual cues.

1 Introduction

005

011

015

017

019

021

042

Humans constantly link matching visual cues to navigate and understand their environment. For instance, we can determine whether objects, individuals are the same simply by comparing their distinguishing visual features (Bruce and Young, 1986; Palermo and Rhodes, 2007; Treisman and Gelade, 1980). This ability, often without needing additional background knowledge, is fundamental in our daily interactions with the world around us. However, while current vision-language models (VLMs) (Chen et al., 2024b; Li et al., 2024b; Zhang



Figure 1: **Previous benchmarks** fail to assess the ability to link matching visual cues, whereas our **VLM²-Bench** explicitly tests this ability, as shown in the example where the model need to identify the reappearance of the same person by linking visual cues, like facial features or clothing, across non-adjacent frames.

et al., 2024b; Team, 2025) have demonstrated extensive knowledge and expanded their capabilities from single-image understanding to handling multiple images and videos, whether thay can effectively link matching visual cues across images or frames—an essential skill for coherent multimodal reasoning—remains an open question.

As shown in Figure 1, existing benchmarks on multiple images and videos fall short in exploring this fundamental ability as they: (a) do not require explicitly linking visual cues across images or frames (Liu et al., 2024c; Yu et al., 2019); (b) rely on external knowledge rather than assessing models' ability to link explicitly visual cues (Zhao et al., 2024; Liu et al., 2024a); (c) emphasize broad and abstract visual comparisons rather than specific cue matching (Wu et al., 2025; Liu et al., 2024b); and (d) focus on retrieval-based tasks rather than eval-



Figure 2: Overview of **VLM**²-**Bench**. The benchmark is categorized into three subsets based on visual cues: GC (General Cue), OC (Object-centric Cue), and PC (Person-centric Cue), each comprising multiple subtasks. To comprehensively evaluate VLMs' ability to visually link matching cues, the benchmark includes diverse question formats—T/F \checkmark , multiple-choice $\frac{1}{60}$, numerical \checkmark , and open-ended ()—ensuring a comprehensive evaluation.

uating the direct association of visual cues across different visual contexts (Wang et al., 2024a).

061

062

063

067

071

087

To bridge this gap, we introduce VLM²-Bench, a benchmark specifically designed to evaluate how well VLMs visually link matching cues. VLM²-Bench is structured around three types of visual cue connection: *general cue*, *person-centric cue*, and *object-centric cue*, encompassing a total of eight subtasks. To balance scalability and quality, we design a semi-automated pipeline with human verification for further refinement. Additionally, our subtasks cover a variety of QA formats—including T/F, multi-choice, numerical, and open-ended questions—totaling over 3,000 question-answer pairs. To better evaluate model performance, we also design specific metrics tailored to various task.

We conduct a comprehensive evaluation of 8 open-source models and GPT-40 on our VLM²-Bench. Despite VLMs generally possessing extensive knowledge, some models perform on par with, or even worse than, the chance-level baseline on our vision-centric tasks. Notably, GPT-40 also underperforms, lagging behind human-level accuracy by 34.80%. This highlights the significant room for improvement in VLMs' ability to link visual cues. Furthermore, we introduce various language-side and vision-side prompting techniques to explore whether they can enhance the models' performance on the benchmark. Through experimental results and case studies, we present *eight key observations*, hoping that these insights will guide future improvements in VLMs for vision-centric tasks.

093

095

099

100

101

103

104

106

107

108

109

110

111

112

113

114

115

116

117

118

119

2 VLM²-Bench

VLM²-Bench is a benchmark designed to assess models' ability to visually link matching cues when processing multiple images or videos. This section introduces the three main categories of VLM²-Bench—general cue (§2.1), object-centric cue (§2.2), and person-centric cue (§2.3)—detailing their associated subtasks, data collection process, and QA pair construction.

2.1 General Cue (GC)

GC is designed to assess a model's ability to link matching cues across diverse contexts, encompassing a broad range of *general cues*. Given two images containing both matched and mismatched cues, an ideal model should accurately identify mismatched ones and associate matched ones.

Subtasks. Here we introduce two subtasks: (i) *Matching (Mat)* evaluates a model's ability to link corresponding visual cues across two images to determine whether they match. Instead of merely identifying differences, the model must associate identical visual elements in both images to recognize what has remained the same and what has changed. (ii) *Tracking (Trk)* focuses on a model's ability to track a specific visual cue that appears in only one of the two images and determine how it has changed. Rather than simply detecting a differ-



Figure 3: Construction of **GC**: (i) We start by manually verifying the edited image data based on three key criteria. (ii) A VLM is then prompted to generate captions for each image, followed by salient score-based filtering to retain the challenging cases. (iii) Finally, visual cues are extracted from two sources and incorporated into a QA prompt, guiding an LLM to generate both positive and negative answer pairs.

ence, the model must link the cue across contexts to understand the transformation process.

Data Collection. We repurpose data from two image editing datasets (Wei et al., 2024; Ku et al., 2023), where each data sample includes an original image I_{ori} , an edited image with subtle modifications I_{edit} , and a corresponding edit instruction \mathcal{P} describing the changes. Our data collection is carried out across two dimensions. First, to ensure diversity in the mismatched cues, GC encompasses various types of changes, such as instance-level modifications (e.g., add/remove, swap, attribute change), which focus on specific items, as well as environment-level changes.

QA Construction. We predefine a T/F question template for *Mat* and *Trk* with a placeholder for the candidate answer (refer to Appendix E). Figure 3 illustrates the construction process, which follows a three-stage approach.

Manual Screening & Refinement: We ensure that \mathcal{P} accurately reflects the changes (correctness), corresponds uniquely to the modified cues (uniqueness), and is unambiguous (clarity).

Salient Sampling: Here, we automate the removal of overly simple cases (e.g., mismatched cues are too salient). To achieve this, a VLM first generates separate descriptions for I_{ori} and I_{edit} , denoted as Cap_{ori} and Cap_{edit} . These descriptions are then combined with \mathcal{P} into a single passage using a predefined template \mathcal{T} (see Table 6 for details). The probability assigned by a language model (e.g., Llama3-8B (Dubey et al., 2024)) to \mathcal{P} given this text-based information is used to compute the salient score, formulated as:

$$S_{\text{salient}} = \frac{1}{|\mathcal{P}|} \sum_{i=1}^{|\mathcal{P}|} \log P_{\theta}(p_i \mid C \cup p_{< i}), \quad (1)$$

where $\mathcal{P} = \{p_1, p_2, ..., p_{|\mathcal{P}|}\}$ represents the tokenized \mathcal{P} , and $C = \mathcal{T}(Cap_{ori}, Cap_{edit})$ denotes the context filled with template \mathcal{T} . Samples with scores below θ (-2.0 here) are retained, ensuring that the benchmark includes more challenging examples requiring nuanced visual cue association.

Pair-wise Answer Generation: Finally, we extract visual cues using a dual-level approach. First, cues parsed from VLM-generated descriptions compensate for the limitations of open-set detectors when handling out-of-distribution scenes. Mean-while, the open-set detector (Wu et al., 2022) extracts fine-grained cues that VLMs might overlook. With these extracted cues, we prompt an LLM to generate a pair of answers for *Mat* and *Trk*, each consisting of one positive and one negative answer.

2.2 Object-centric Cue (OC)

OC aims to assess a model's ability to link matching cues associated with everyday objects using *object-centric cues*. Even when encountering an object for the first time, a well-aligned model should be able to leverage its unique visual cues to establish associations, enabling it to recognize and track the object across different scenes. This capability is essential for coherent perception and interaction in real-world deployments.

Subtasks. Based on the complexity of linking cues to solve the problem, we define three subtasks in OC. (i) *Comparison (Cpr)* requires the model to determine whether the objects appearing in different images are the same. This task

primarily assesses the model's ability to perceive 186 visual consistency or change. Notably, we observe 187 that models exhibit significant model-specific bias when making a binary decision (Goyal et al., 2017; Ye et al., 2024b; Song et al., 2024; Li et al., 2024a), leading to discrepancies between results and their 191 actual capabilities. To mitigate this, we introduce 192 consistency-pair validation, where for each state-193 ment (e.g., "X is Y", with the answer being T), 194 we generate a corresponding negation (e.g., "X is 195 not Y", with the answer being F). The model is only considered correct if it correctly answers both 197 statements, ensuring consistency in its decision-198 making. (ii) *Counting (Cnt)* involves identifying 199 the number of unique objects, requiring the model 200 not only to recognize variations or consistencies but also to track distinct cues to avoid double-counting the same object. (iii) Grouping (Grp), the most challenging one, requires the model to identify all 204 instances of the same object, building on precise 205 cue matching across multiple images.

Data Collection. We manually collect various categories of everyday objects (e.g., pets, cups). For each category, we define multiple subcategories and collect a set of images \mathcal{I}_{O_i} —four images that depict the same object in different scenarios. Additionally, we also collect a set $\mathcal{I}_{\neg O_i}$, consisting of four images of different objects, each containing some matching visual cues with \mathcal{I}_{O_i} , which are used as distractors.

> **QA Construction.** For each subtask, we define a question template that includes a placeholder for \mathcal{I}_{O_i} , which allows us to tailor the question based on different objects (see Appendix E). For answer generation, we first curate the multi-image sequences according to predefined rules. For each specific sequence, we generate the ground truth answers for the questions related to *Cpr*, *Cnt*, and *Grp*.

2.3 Person-centric Cue (PC)

216

217

218

219

220

221

223

224

225

233

PC aims to evaluate a model's ability to link *personcentric cues*. While a model cannot memorize every individual, it should possess the capability to associate the same person across different images or frames by leveraging distinctive visual cues such as facial features, clothing, or body posture. This ability is essential for ensuring coherent perception of human actions and is a fundamental requirement for real-world VLM applications.

Subtasks. Similar to OC's subtasks (refer to §2.2), PC includes (i) *Comparison (Cpr)*, (ii) *Counting (Cnt)*, and (iii) *Grouping (Grp)*. However, unlike objects, individuals can be observed through their actions in videos. Therefore, we introduce (iv) *Video Identity Describing (VID)*. This subtask assesses whether a model can correctly link the same person by analyzing its description of a video containing that person.

Data Collection. We manually select several individuals, each denoted as \mathcal{P}_i . For each individual, we collect $\mathcal{I}_{\mathcal{P}_i}$ —4 images depicting the same individual. For each image $I_i \in \mathcal{I}_{\mathcal{P}_i}$, we select the distractor images $I_{\neg i} \notin \mathcal{I}_{\mathcal{P}_i}$ that has the highest CLIP similarity (Hessel et al., 2021). This allows us to obtain images of different individuals where most cues are matched. For the subtask of VID, we collect videos of different individuals, denoted as $V_{\mathcal{P}_i}$, and pair each with another video $V_{\neg \mathcal{P}_i}$ featuring a different individual with highly similar cues (e.g., actions, scene, clothing). We then construct two video sequences: (i) $\mathcal{P}_i \rightarrow \neg \mathcal{P}_i$, assessing the model's ability to distinguish individuals. (ii) $\mathcal{P}_i \rightarrow \neg \mathcal{P}_i \rightarrow \mathcal{P}_i$, evaluating whether the model detects changes and links the final occurrence of \mathcal{P}_i to its first appearance.

QA Construction. The construction for the overall QA in PC's *Cpr*, *Cnt*, and *Grp* subtasks follows a similar approach to OC. For the *VID* task, we emphasize the model's ability to describe individuals when designing open-ended questions, aiming to better test the model's capacity to link individuals appearing in different scenes.



Figure 4: Statistical overview of **VLM²-Bench**. The pie chart shows the distribution of 9 subtasks across the 3 main categories of visual cues. The bar plot illustrates the percentage breakdown by question format.

2.4 Benchmark Statistics

Our benchmark is organized into three main categories, comprising a total of 9 subtasks. After careful verification, it contains 3,060 question-answer 267 268 269

270

234

235

236

237

238

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

264

265

Baselines or Models	G	C		OC			I	PC		Ov	erall*
	Mat	Trk	Cpr	Cnt	Grp	Cpr	Cnt	Grp	VID	Avg	Δ_{human}
Chance-Level	25.00	25.00	50.00	34.88	25.00	50.00	34.87	25.00	-	33.72	-61.44
Human-Level	95.06	98.11	96.02	94.23	91.92	97.08	92.87	91.17	100.00	95.16	0.00
LLaVA-OneVision-7B	16.60	13.70	47.22	56.17	27.50	62.00	46.67	37.00	47.25	39.35	-55.81
LLaVA-Video-7B	18.53	12.79	54.72	62.47	28.50	62.00	66.91	25.00	59.00	43.32	-51.84
LongVA-7B	14.29	19.18	26.67	42.53	18.50	21.50	38.90	18.00	3.75	22.59	-72.57
mPLUG-Owl3-7B	17.37	18.26	49.17	62.97	31.00	63.50	58.86	26.00	13.50	37.85	-57.31
Qwen2-VL-7B	27.80	19.18	68.06	45.99	35.00	61.50	58.59	49.00	16.25	42.37	-52.79
Qwen2.5-VL-7B	35.91	43.38	71.39	41.72	47.50	80.00	57.98	69.00	46.50	54.82	-40.34
InternVL2.5-8B	21.24	26.03	53.33	55.23	46.50	51.50	60.00	52.00	5.25	41.23	-53.93
InternVL2.5-26B	30.50	30.59	43.33	51.48	52.50	59.50	59.70	61.00	21.75	45.59	-49.57
GPT-40	37.45	39.27	74.17	80.62	57.50	50.00	90.50	47.00	66.75	60.36	-34.80

Table 1: Evaluation results on VLM²-Bench, covering *Mat* (Matching), *Trk* (Tracking), *Cpr* (Comparison), *Cnt* (Counting), *Grp* (Grouping), and *VID* (Video Identity Describing). The highest, second, and third highest scores are highlighted. *: Overall excludes the *VID* due to the lack of a chance-level baseline for open-ended tasks.

pairs, with varying formats including T/F, multichoice (MC), numerical (Nu), and open-ended (Oe). To ensure the quality of the annotations, we perform an inter-annotator agreement (IAA) evaluation (Thorne et al., 2018) involving three annotators, resulting in a high Fleiss' Kappa score (Fleiss, 1971) of 0.983. Figure 4 presents the distribution of these subtasks across the three categories, along with the breakdown of different question formats. For additional details, refer to Appendix C.

3 Evaluation

271

272

273

275

276

277

278

279

283

287

290

291

292

294

296

3.1 Metric Design

T/F (*Matching, Tracking, Comparison*): Accuracy is computed based on paired evaluation, where a response is correct only if it answers T (ground-truth True) and F (ground-truth False) correctly. The overall accuracy across N test pairs is:

$$Acc_{pair} = \frac{\sum_{i=1}^{N} \left(T_i^+ \cap F_i^-\right)}{N}, \qquad (2)$$

where T^+ and F^- denote correct predictions for T and F, respectively.

Numerical (*Counting*): Absolute matching alone does not effectively reflect the severity of errors in numerical responses. To measure the extent of the error between the predicted count \hat{N}_i and ground truth N_i , we introduce Acc_{num} . The first step is to calculate the normalized error:

$$\epsilon_i = \frac{\left|\hat{N}_i - N_i\right|}{\max\left(N_i - 1, N_i^{img} - N_i\right)},\qquad(3)$$

where N_i^{img} is the number of input images. We define $w_i = \max(\{N_i^{img}\}_{i=1}^n)/N_i^{img}$ to penalize errors in cases with fewer images and introduce α as an error amplification factor. The final accuracy over n cases is:

$$Acc_{num} = 1 - \frac{1}{n} \sum_{i=1}^{n} w_i \cdot \epsilon_i^{-\alpha}.$$
 (4)

298

299

300

301

302

303

304

305

306

307

308

310

311

312

313

314

315

316

317

318

319

321

322

323

324

325

327

Multi-choice (*Grouping*): Accuracy is the proportion of correctly predicted choices.

Open-ended (*Video Identity Describing*): We use GPT-40 to score model's descriptions, in combination with rule-based scoring prompts. The final accuracy Acc_{oe} is obtained by averaging the scores of all open-ended responses and rescaling them to the range of [0,1]. Additionally, we perform manual verification of GPT-40's scoring. For each model, we randomly sample 20 scored responses for review, and find only 2 instances with discrepancies, resulting in an accuracy rate of 98.89% (178/180). Refer to Appendix F for more details.

3.2 Evaluation Setup

Evaluated Models. We evaluate eight opensource VLMs that support multiple-image or video input: LLaVA-OneVision (Li et al., 2024b), LLaVA-Video (Zhang et al., 2024b), LongVA (Zhang et al., 2024a), mPLUG-Owl3 (Ye et al., 2024a), Qwen2-VL (Wang et al., 2024b), Qwen2.5-VL (Team, 2025), and InternVL2.5 (Chen et al., 2024b). Additionally, we include the commercial model GPT-40 (Hurst et al., 2024) for comparison.

Model		Matching (Mat) Tracking (Trk)						
	A/R	Swp	Attr	Env	A/R	Swp	Attr	Env
LV-OV	50.68	49.15	53.45	52.50	27.27	45.51	57.50	70.59
LV-Vid	56.08	49.15	53.45	51.25	46.75	48.88	52.50	67.65
LongVA	37.84	46.58	53.45	46.25	46.10	49.44	42.50	60.29
Owl3	54.73	52.56	55.17	50.00	41.56	48.88	55.00	73.53
Qw2-VL	53.68	52.56	55.17	68.75	65.58	62.90	77.50	63.93
Qw2.5-VL	64.19	55.62	74.14	67.50	61.69	69.10	55.00	64.71
In2.5-8B	64.86	51.28	52.07	66.25	54.55	67.42	62.50	60.65
In2.5-26B	60.81	51.71	58.62	61.25	56.49	62.92	47.50	66.18
GPT-40	75.00	61.97	56.90	70.00	68.83	67.98	67.50	64.71

Table 2: Breakdown of four mis-matched cue types in two subtasks of GC. For each model, the highest and second highest error (%) per subtask are highlighted.

Baselines. We introduce chance-level and human-level baselines (details are in Appendix D).

3.3 Results and Findings

328

329

330

331

335

338

340

342

343

345

347

Results. Table 1 presents the comprehensive performance of various models across the three categories - General Cue (GC), Object-centric Cue (OC), and Person-centric Cue (PC) – of our VLM²-Bench, covering a total of nine subtasks.

Finding I: Simple tasks for humans pose sig-336 nificant challenges for VLMs. We observe that humans achieve near-perfect accuracy across most tasks in our VLM²-Bench. In contrast, even GPT-40, a state-of-the-art model, performs significantly lower than humans, with an overall performance gap of 34.80%. For open-source models, many show performance comparable to the chance-level baseline or only slightly outperform it. Specifically, for the VID, humans can easily achieve 100% accuracy in distinguishing and linking individuals in a video. However, even the best-performing model, GPT-40, reaches only 66.75%. Errors mainly arise from failing to recognize individuals after changes or misidentifying reappearing persons as new.

Finding II: Relatively consistent error patterns in Mat and Trk of GC. Table 2 shows that models struggle with mismatched cues due to swap in Mat, which requires linking two completely different cues. To identify what has changed, models must first link and match all the other cues in the 356 context before they can determine that the swapped cue has been transformed. This task requires a deeper understanding of how cues relate to each other across different instances. In contrast, Trk challenges models with mismatched cues due to ad-361 d/remove, which focuses on tracking how a specific cue changes. This suggests that when there is a cue that appears only once, the model struggles to link 364

the non-appearing cue with the appearing cue to track the transformation process effectively. This limitation reveals models' difficulty in handling cases where certain cues are missing but still need to be linked to understand the dynamic changes.

365

366

367

369

370

371

372

373

374

375

376

377

379

380

381

385

386

387

388

391

392

393

394

395

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

Finding III: Models perform better in linking person-centric cues than object-centric cues. We selected the top three open-source models (Qwen2.5-VL-8B, InternVL2.5-8B, InternVL2.5-26B) and compared their performance on the three shared tasks (Cpr, Cnt, Grp) in both OC and PC. Results show that, on average, the performance on PC is higher than on OC by 7.65%, 9.75%, and 11.83% for the tasks of Cpr, Cnt, Grp, respectively. This could be due to the fact that, during training on person-related data, models are likely provided with explicit person names as anchors to person-centric cues, which helps the models better distinguish different individuals. In contrast, objects are typically trained using general category names, which may not provide such clear distinctions. Additionally, these models might have been specifically trained on large datasets that emphasize differentiating and linking individuals (Pi et al., 2024a; Dai et al., 2024), thereby enhancing their ability to link person-centric cues.

How Prompting Methods affect VLMs 4

In this section¹, we investigate various prompting methods (language-side and vision-side) to evaluate their impact on performance in VLM²-Bench. We select the top 3 performing open-source models (Qwen2.5-VL-8B, InternVL2.5-8B, InternVL2.5-26B), along with GPT-40, and explore different approaches of CoT (Kojima et al., 2022; Wei et al., 2023) and visual prompting (VP) (Lei et al., 2024; Yang et al., 2023) (refer to Appendix F for details). The goal is to investigate whether these techniques can improve performance across the benchmark and to identify the underlying factors that contribute to their success or failure.

4.1 Probing for General Cue (GC)

Methods. (i) CoT-normal (Table 23) encourages the model to solve the task step by step, allowing it to reason through the problem. (ii) CoT-special (Table 24) guides the model to solve the task using a thought process closer to how humans typically approach it. (iii) VP-grid (Figure 11) is adapted

¹Due to space limits, we reference most case studies, figures, and details in the Appendix within this section.



InternVL2.5-8B

VP-zoom-p

InternVL2.5-26B

🚥 Cnt

GPT-40

Figure 5: Performance gains or losses (%) when applying different prompting methods on VLM²-Bench.

from previous work (Lei et al., 2024) for our tasks, 412 overlaying a dot matrix on the image as visual an-413 chors to provide positional references and enhance 414 the model's performance. 415

-30

416

417

418

419

420

421

422

423

424

425

426

427

428

429

en2.5-VL-7B

Finding IV: Reasoning in language aids models in logically linking visual cues. From Figure 5a, it is evident that both CoT-normal and CoT-special, which reasoning in language, positively impact model performance in most cases. As demonstrated in Figure 14, CoT-special improves performance by first having the model explicitly write out the cues present in each image, followed by using language to make inferences. This process helps reduce the model's error rate by structuring the task and providing clearer logical guidance. This suggests that when models are linking general visual cues, using language to help structure the logical flow of the process can be beneficial.

Finding V: Effectiveness of visual prompting 430 depends on models' ability to interpret both 431 prompting cues and the visual content. As 432 shown in Figure 5a, VP-grid negatively impacts GC 433 performance for QwenVL2.5, causing a significant 434 435 drop compared to the vanilla approach. Figure 15 reveals that this decline stems from the model's dif-436 ficulty in interpreting the visual coordinates within 437 the prompt, leading to misinterpretation of the cues 438 and causing it to fail cases it originally answered 439

correctly under the vanilla setting. However, as shown in Figure 16, GPT-40 successfully resolves a previously incorrect case by effectively leveraging the cues introduced through visual prompting while utilizing its strong visual perception abilities.

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

4.2 Probing for Object-centric Cue (OC)

Methods. (i) CoT (Table 23). (ii) VP-zoom-o (Figure 12) uses an open-set detector (Ren et al., 2024) to obtain bounding boxes, which are then cropped to focus the model's attention on objectcentric cues. By eliminating irrelevant non-object cues and emphasizing the object-centric cues, it enhances the model's ability to better focus on the most relevant visual information.

Finding VI: The open-ended nature of language may hinder object grouping. Unlike GC that link instance-level cues, OC requires grouping similar objects based on fine-grained visual details. As shown in Figure 5b, InternVL2.5 using CoT struggles with this task because the open-ended nature of language leads to both limited coverage of subtle visual cues (see Figure 17) and inconsistent representations of the same cues, introducing ambiguity, making it harder for models to reliably align and group matching objects.

Finding VII: Amplifying object cues benefits stronger models while having minimal impact on others. From Figure 5b, we observe that for models with strong vision capabilities like GPT-40, our VP-zoom-o method further enhances performance. For other models, this method at least ensures that the performance remains on par with the vanilla approach, without causing any degradation.

4.3 Probing for Person-centric Cue (PC)

Methods. (i) CoT (Table 23). (ii) VP-zoom-p (Figure 13) utilizes a face detector (Geitgey, 2016) to obtain bounding boxes of faces-the most distinguishing feature of different individuals. It then crops the image to focus only on the face, thereby minimizing the interference from distractor cues such as clothing and other background elements.

Finding VIII: CoT and visual prompting fail to improve linking on highly abstract personcentric cues, leading to a performance drop. From Figure 5c, we observe that for almost all models, neither CoT (language-based) nor VP-zoom-p (vision-based) lead to improved performance. This is because facial features are highly abstract, and

⁽c) Results of CoT and VP-zoom-p on PC.

579

580

581

582

583

584

585

538

539

CoT methods struggle to effectively describe them in words. Additionally, VP-zoom-p fails because current models' visual capabilities are insufficient to accurately perceive facial features.

5 Related Work

488

489

490

491

492

493

494

495

497

498

499

503

504

506

507

508

509

510

511

512

513

515

516

517

518

Recent advancements in vision-language models (Hurst et al., 2024; Team, 2025; Zhang et al., 2024a; Li et al., 2024b; Ye et al., 2024a; Chen et al., 2024b; Liang et al., 2024b) have significantly broadened their capabilities. Previously restricted to processing single-image inputs, many VLMs can now handle multi-image and even video inputs, allowing them to capture richer and more dynamic visual contexts. Additionally, with access to a growing volume of high-quality visual-textual paired training data (Pi et al., 2024b; Garg et al., 2024; Chen et al., 2023; Zhang et al., 2024c; Wang et al., 2024c), these models have shown substantial improvements in perceiving subtle visual cues and their relationships, enabling them to engage in more nuanced reasoning about visual content. Furthermore, VLMs are increasingly applied in real-world scenarios, including navigation (Weerakoon et al., 2024), planning (Yang et al., 2024), and autonomous driving (Jiang et al., 2024), solidifying their role in bridging vision and language for practical applications. However, to truly integrate into everyday life, VLMs still have significant room for improvement when it comes to more fundamental but common visual tasks, such as those assessed in our benchmark.

Benchmarking vision-language models plays 519 a critical role in guiding their future develop-520 ment (Liang et al., 2024a; Yin et al., 2023; Chen 521 et al., 2024a). These benchmarks typically focus on 522 assessing the models' fine-grained perception (Li et al., 2024a; Tong et al., 2024), reasoning abil-524 ities (Lu et al., 2022; Yu et al., 2023; Wu et al., 525 2024), commonsense knowledge (Yue et al., 2024). 526 In addition, evaluations targeting multi-image and 527 video inputs are designed to measure the new competencies that VLMs require as their visual context extends. These tasks include captioning (Yue et al., 2024; Yu et al., 2019), retrieval (Wang et al., 2024a; 531 Li et al., 2025), comparison (Wu et al., 2025; Jiao 533 et al., 2024), and temporal reasoning (Liu et al., 2024b). However, existing benchmarks focus on 534 evaluating VLMs' ability to interpret visual cues based on their knowledge. In contrast, humans typically solve such tasks by explicitly matching vi-537

sual cues without relying on extensive background knowledge. To better assess whether they can replicate this human-like ability, we propose VLM²-Bench, which focuses on linking and matching explicit visual cues.

6 Takeaways

Based on our findings, we highlight three key areas for future improvements:

- Strengthening Fundamental Visual Capabilities. Improving core visual abilities not only enhances overall performance but also increases adaptability. A stronger visual foundation maximizes the effectiveness of visual prompting and reduces reliance on prior knowledge, enabling models to operate more independently in vision-centric tasks.
- Balancing Language-Based Reasoning in Vision-Centric Tasks. Integrating language into vision-centric tasks requires careful calibration. Future research should establish clearer principles on when language-based reasoning aids visual understanding and when it introduces unnecessary biases, ensuring models leverage language appropriately.
- Evolving Vision-Text Training Paradigms. Current training paradigms focus heavily on emphasizing vision-language associations. However, as models expand their visual context window, their ability to reason purely within the visual domain becomes increasingly crucial. We should prioritize developing models that can structure, organize, and infer relationships among visual cues.

7 Conclusion

In summary, we introduce VLM²-Bench, a novel benchmark designed to probe the capability of vision-language models (VLMs) in visually linking matching cues, an essential yet underexplored skill for models in everyday visual reasoning. Through extensive evaluations and further analysis of prompting techniques applied on our benchmark, we identify 8 key findings. Notably, even GPT-40 falls 34.80% behind human performance. Based on these insights, we advocate for advancements in fundamental visual capabilities, better integration of language-based reasoning, and the evolution of vision-text training paradigms to improve VLMs' performance in vision-centric tasks.

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

Limitations

586

599

605

606

607

610

611

612

613

614

615

616

617

618

619

623

625

627

630

631

632

636

587VLM2-Bench focuses on evaluating visual cue link-588ing but does not cover all possible scenarios. Ad-589ditionally, while it provides valuable insights, its590scale is limited, and model performance may not591fully generalize to all real-world settings. Auto-592mated evaluation constraints limit the inclusion of593open-ended questions in our benchmark, impacting594the assessment of models' vision-centric reason-595ing abilities. Expanding task diversity and refining596evaluation methods remain important directions for597future work.

References

- Vicki Bruce and Andrew W Young. 1986. Understanding face recognition. *British journal of psychology*, 77 (Pt 3):305–27.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. 2023. Sharegpt4v: Improving large multimodal models with better captions. *Preprint*, arXiv:2311.12793.
- Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. 2024a. Are we on the right way for evaluating large vision-language models? *arXiv preprint arXiv:2403.20330*.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 2024b. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24185–24198.
- Dawei Dai, Xu Long, Li Yutang, Zhang Yuanhui, and Shuyin Xia. 2024. Humanvlm: Foundation for human-scene vision-language model. *Preprint*, arXiv:2411.03034.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378.
- Roopal Garg, Andrea Burns, Burcu Karagol Ayan, Yonatan Bitton, Ceslee Montgomery, Yasumasa Onoe, Andrew Bunner, Ranjay Krishna, Jason Baldridge, and Radu Soricut. 2024. Imageinwords: Unlocking hyper-detailed image descriptions. *Preprint*, arXiv:2405.02793.

- Adam Geitgey. 2016. Machine learning is fun! part
 4: Modern face recognition with deep learning. *Medium. Medium Corporation*, 24:2016.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. Clipscore: A referencefree evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- 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. *Preprint*, arXiv:2410.22313.
- Qirui Jiao, Daoyuan Chen, Yilun Huang, Yaliang Li, and Ying Shen. 2024. Img-diff: Contrastive data synthesis for multimodal large language models. *arXiv preprint arXiv:2408.04594*.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199– 22213.
- Max Ku, Tianle Li, Kai Zhang, Yujie Lu, Xingyu Fu, Wenwen Zhuang, and Wenhu Chen. 2023. Imagenhub: Standardizing the evaluation of conditional image generation models. *arXiv preprint arXiv:2310.01596*.
- Xuanyu Lei, Zonghan Yang, Xinrui Chen, Peng Li, and Yang Liu. 2024. Scaffolding coordinates to promote vision-language coordination in large multi-modal models. *Preprint*, arXiv:2402.12058.
- Baiqi Li, Zhiqiu Lin, Wenxuan Peng, Jean de Dieu Nyandwi, Daniel Jiang, Zixian Ma, Simran Khanuja, Ranjay Krishna, Graham Neubig, and Deva Ramanan. 2024a. Naturalbench: Evaluating visionlanguage models on natural adversarial samples. *arXiv preprint arXiv:2410.14669*.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. 2024b. Llavaonevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- You Li, Heyu Huang, Chi Chen, Kaiyu Huang, Chao Huang, Zonghao Guo, Zhiyuan Liu, Jinan Xu, Yuhua Li, Ruixuan Li, et al. 2025. Migician: Revealing

multimodal large language models. arXiv preprint Large language models are skeptics: False negative 746 arXiv:2501.05767. problem of input-conflicting hallucination. arXiv 747 preprint arXiv:2406.13929. 748 Paul Pu Liang, Akshay Goindani, Talha Chafekar, Leena Mathur, Haofei Yu, Ruslan Salakhutdinov, and Louis-Qwen Team. 2025. Qwen2.5-vl. 749 Philippe Morency. 2024a. Hemm: Holistic evaluation of multimodal foundation models. arXiv Thorne, Vlachos, James Andreas Christos 750 preprint arXiv:2407.03418. Christodoulopoulos, and Arpit Mittal. 2018. 751 Fever: a large-scale dataset for fact extraction and 752 Paul Pu Liang, Amir Zadeh, and Louis-Philippe verification. arXiv preprint arXiv:1803.05355. 753 Morency. 2024b. Foundations & trends in multimodal machine learning: Principles, challenges, and Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, 754 open questions. ACM Computing Surveys, 56(10):1-Yann LeCun, and Saining Xie. 2024. Eyes wide 755 shut? exploring the visual shortcomings of multi-756 modal llms. In Proceedings of the IEEE/CVF Con-757 ference on Computer Vision and Pattern Recognition, 758 pages 9568-9578. 759 Anne Treisman and Garry A. Gelade. 1980. A feature-760 integration theory of attention. Cognitive Psychology, 761 12:97-136. 762 Fei Wang, Xingyu Fu, James Y Huang, Zekun Li, Qin 763 Liu, Xiaogeng Liu, Mingyu Derek Ma, Nan Xu, 764 arXiv preprint Wenxuan Zhou, Kai Zhang, et al. 2024a. Muirbench: 765 A comprehensive benchmark for robust multi-image 766 understanding. arXiv preprint arXiv:2406.09411. 767 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-768 hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin 769 Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024b. Qwen2-vl: Enhancing vision-language model's per-773 ception of the world at any resolution. arXiv preprint 774 arXiv:2409.12191. 775 Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo 776 Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan Chen, 777 Yaohui Wang, Conghui He, Ping Luo, Ziwei Liu, 778 Yali Wang, Limin Wang, and Yu Qiao. 2024c. In-779 ternvid: A large-scale video-text dataset for mul-780 timodal understanding and generation. Preprint, 781 arXiv:2307.06942. 782 Kasun Weerakoon, Mohamed Elnoor, Gershom Senevi-783 ratne, Vignesh Rajagopal, Senthil Hariharan Arul, 784 Jing Liang, Mohamed Khalid M Jaffar, and Dinesh Manocha. 2024. Behav: Behavioral rule guided autonomy using vlms for robot navigation in outdoor 787 scenes. Preprint, arXiv:2409.16484. Cong Wei, Zheyang Xiong, Weiming Ren, Xinrun Du, 789 Ge Zhang, and Wenhu Chen. 2024. Omniedit: Build-790 ing image editing generalist models through special-791 ist supervision. arXiv preprint arXiv:2411.07199. 792 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten 793 Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and 794 Denny Zhou. 2023. Chain-of-thought prompting elic-795 its reasoning in large language models. Preprint, 796 arXiv:2201.11903. 797

Jongyoon Song, Sangwon Yu, and Sungroh Yoon. 2024.

745

42.

the magic of free-form multi-image grounding in

697

699

700

701

703

704

705

707

710

713

714

715

716

717

718

719

720

721 722

723

724

726

727

730

731

733

736

737

738

739

740

741

742

743

744

- Haowei Liu, Xi Zhang, Haiyang Xu, Yaya Shi, Chaoya Jiang, Ming Yan, Ji Zhang, Fei Huang, Chunfeng Yuan, Bing Li, et al. 2024a. Mibench: Evaluating multimodal large language models over multiple images. arXiv preprint arXiv:2407.15272.
- Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. 2024b. Tempcompass: Do video llms really understand videos? arXiv:2403.00476.
- Ziyu Liu, Tao Chu, Yuhang Zang, Xilin Wei, Xiaoyi Dong, Pan Zhang, Zijian Liang, Yuanjun Xiong, Yu Qiao, Dahua Lin, et al. 2024c. Mmdu: A multi-turn multi-image dialog understanding benchmark and instruction-tuning dataset for lvlms. arXiv preprint arXiv:2406.11833.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. Advances in Neural Information Processing Systems, 35:2507–2521.
- Romina Palermo and Gillian Rhodes. 2007. Are you always on my mind? a review of how face perception and attention interact. Neuropsychologia, 45:75–92.
- Renjie Pi, Jianshu Zhang, Tianyang Han, Jipeng Zhang, Rui Pan, and Tong Zhang. 2024a. Personalized visual instruction tuning. arXiv preprint arXiv:2410.07113.
- Renjie Pi, Jianshu Zhang, Jipeng Zhang, Rui Pan, Zhekai Chen, and Tong Zhang. 2024b. Image textualization: An automatic framework for creating accurate and detailed image descriptions. arXiv preprint arXiv:2406.07502.
- Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang Li, Qing Jiang, and Lei Zhang. 2024. Grounded sam: Assembling openworld models for diverse visual tasks. Preprint, arXiv:2401.14159.

853

873

- 798 799

European Conference on Computer Vision, pages

Jialian Wu, Jianfeng Wang, Zhengyuan Yang, Zhe

Gan, Zicheng Liu, Junsong Yuan, and Lijuan Wang.

former for object understanding. arXiv preprint

Shujin Wu, Yi Fung, Sha Li, Yixin Wan, Kai-Wei Chang,

and Heng Ji. 2024. MACAROON: Training vision-

language models to be your engaged partners. In

Findings of the Association for Computational Lin-

guistics: EMNLP 2024, pages 7715-7731, Miami,

Florida, USA. Association for Computational Lin-

Lingfeng Yang, Yueze Wang, Xiang Li, Xinlong Wang,

Zhutian Yang, Caelan Garrett, Dieter Fox, Tomás

Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou.

2024a. mplug-owl3: Towards long image-sequence understanding in multi-modal large language models.

Jiavi Ye, Yanbo Wang, Yue Huang, Dongping Chen, Qihui Zhang, Nuno Moniz, Tian Gao, Werner Geyer,

Chao Huang, Pin-Yu Chen, et al. 2024b. Justice

or prejudice? quantifying biases in llm-as-a-judge.

Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing

Sun, Tong Xu, and Enhong Chen. 2023. A survey on

multimodal large language models. arXiv preprint

Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. arXiv preprint

Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. 2019. Activitynet-qa: A dataset for understanding complex web videos via question answering. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 33, pages

Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2024. Mmmu: A massive multi-discipline multimodal understanding

ning with vision language models.

Lozano-Pérez, and Leslie Pack Kaelbling. 2024. Guiding long-horizon task and motion plan-

Preprint,

and Jian Yang. 2023. Fine-grained visual prompting.

Grit: A generative region-to-text trans-

360–377. Springer.

arXiv:2212.00280.

Preprint, arXiv:2306.04356.

Preprint, arXiv:2408.04840.

arXiv preprint arXiv:2410.02736.

arXiv:2410.02193.

arXiv:2306.13549.

arXiv:2308.02490.

9127-9134.

2022.

guistics.

- 809 810 811
- 812
- 814
- 816

817

- 819
- 822
- 825 826
- 827 830

831 832

- 834 835 836 837
- 842

847

848

852

Haoning Wu, Hanwei Zhu, Zicheng Zhang, Erli Zhang, and reasoning benchmark for expert agi. In Proceedings of the IEEE/CVF Conference on Computer Chaofeng Chen, Liang Liao, Chunyi Li, Annan Wang, Wenxiu Sun, Qiong Yan, et al. 2025. To-Vision and Pattern Recognition, pages 9556–9567. wards open-ended visual quality comparison. In

- Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. 2024a. Long context transfer from language to vision. arXiv preprint arXiv:2406.16852.
- Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. 2024b. Video instruction tuning with synthetic data. Preprint, arXiv:2410.02713.
- Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. 2024c. Video instruction tuning with synthetic data. Preprint, arXiv:2410.02713.
- Bingchen Zhao, Yongshuo Zong, Letian Zhang, and Timothy Hospedales. 2024. Benchmarking multiimage understanding in vision and language models: Perception, knowledge, reasoning, and multi-hop reasoning. arXiv preprint arXiv:2406.12742.

874	A Appendix Outline
875	In the appendix, we provide:
876	• Appendix B provides details on the licensing
877	terms and usage rights for our benchmark.
878	• Appendix C presents the statistical analysis
879	of the VLM ² -Bench.
880	• Appendix D details on how we obtain the
881	chance-level and human-level baselines.
882	• Appendix E elaborates more details on the
883	construction of the VLM ² -Bench.
884	• Appendix F provides a deeper dive into the
885	various prompting techniques we use.
886	• Appendix G a detailed breakdown and analy
887	sis of failure and success examples regarding
888	different prompting methods.
889	B Licencing and Intended Use
890	Our VLM ² -Bench is available under the CC-BY
891	4.0 license for academic use with proper attribu-
892	tion. The images, videos, and annotations in this
893	benchmark are intended solely for research pur
894	poses. These data were sourced from publicly avail
895	able online platforms, and while efforts were made
896	to use them responsibly, explicit permissions may
897	not have been obtained for all content. Users are
898	responsible for ensuring that their use of the data
899	complies with applicable intellectual property laws
900	and ethical guidelines. We encourage users to ver
901	ity the sources and ensure compliance with any
902	terms of service or licensing agreements.

VLM²-Bench Statistics С

903

Here we provide additional details regarding the 904 construction and statistics of our VLM²-Bench 905 benchmark. As described in the main paper (§ 2.4), 906 our benchmark comprises three main categories-907 General Cue (GC), Object-centric Cue (OC), and 908 Person-centric Cue (PC)—with a total of 3,060 909 visual-text query pairs. Below, we elaborate on 910 the specific data composition, including the dis-911 tribution of question types (T/F, multiple-choice 912 (MC), numerical (Nu), and open-ended (Oe)) and 913 914 the rationale behind each subtask.

Category	Subtask	T/F	MC	Nu	Oe	Total
CC	Mat	520	_	_	_	520
UC.	Trk	440	-	-	_	440
	Subtotal	960	_	_	-	960
	Cpr	720	_	_	_	720
OC	Cnt	-	-	360	_	360
	Grp	-	200	-	-	200
	Subtotal	720	200	360	-	1,280
	Cpr	400	_	_	_	400
DC	Cnt	_	-	120	_	120
PC	Grp	-	100	-	-	100
	VID	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	200	200		
	Subtotal	400	100	120	200	820
Total		2,480	300	480	200	3,060

Table 3: Overview of query distribution across the three categories of VLM²-Bench. T/F = True/False, MC = multiple-choice, Nu = numerical, Oe = open-ended.

C.1 Overall Composition

Table 3 summarizes the total query counts within each category and subtask. The benchmark is divided as follows:

915

916

917

918

939

• General Cue (GC): 960 queries	919
- <i>Matching (Mat)</i> : 260 T/F pairs \rightarrow 520 queries	920 921
- <i>Tracking (Trk)</i> : 220 T/F pairs \rightarrow 440 queries	922
• Object-centric Cue (OC): 1,280 queries	923
- Comparison (Cpr): 360 T/F pairs \rightarrow 720 queries	924 925
- <i>Counting (Cnt)</i> : 360 numerical (counting) queries	926 927
- <i>Grouping (Grp)</i> : 200 multiple-choice questions	928 929
• Person-centric Cue (PC): 820 queries	930
- Comparison (Cpr): 200 T/F pairs \rightarrow 400 queries	931 932
- <i>Counting (Cnt)</i> : 120 numerical (counting) queries	933 934
- <i>Grouping (Grp)</i> : 100 multiple-choice questions	935 936
- Free-form (VID): 200 open-ended queries	937
Summing these yields a total of 3 060 visual-text	029

Summing these yields a total of 3,060 visual-text query pairs.

990

991

C.2 Details per Subtask and Question Type

General Cue (GC).

941

Matching (Mat). We collect 260 True/False (T/F)
pairs focused on verifying the alignment between
a visual instance and a textual description (e.g.,
object presence, basic attributes). Each T/F pair
forms two distinct queries (one True, one False),
yielding 520 queries in total.

948Tracking (Trk). We design 220 T/F pairs that949test an understanding of object or entity continu-950ity across frames. For example, a question might951ask whether the same object reappears in subse-952quent frames. Each T/F pair similarly results in953two queries, totaling 440.

954Object-centric Cue (OC). All the visual query955cases are built upon the 360 image sequences we956construct. Details about image sequences can be957found in Section E.2.

Comparison (Cpr). This subtask examines the model's ability to compare object properties (e.g., size, color, quantity) across different frames. We produce 360 T/F pairs, each yielding two queries (720 total). Among these 360 pairs, we maintain a 1:2 ratio of True to False for ground-truth answers (i.e., 120 True vs. 240 False).

Counting (Cnt). We provide 360 numerical questions, each asking for a count of objects in a given
scene or sequence. Possible numeric answers are
typically small integers (e.g., 1, 2, 3), reflecting the
number of relevant objects.

Grouping (Grp). We generate 200 multiple-choice 970 (MC) questions that ask about grouping objects ac-971 cording to certain criteria (e.g., AAB, ABC, AAAB, 972 AABC, ABCD). Each question presents multiple 973 group-configuration options plus a "None" option, 974 which can serve as either a correct or distractor 975 choice. For image sequences of length 4, the op-976 tions include various plausible groupings (two-of-a-977 kind, three-of-a-kind, etc.) along with at least one 978 additional distractor grouping that also involves three-of-a-kind to ensure sufficient challenge.

Person-centric Cue (PC). Similar to OC, the construction of 260 image sequences as well as 200 video clips for PC is detailed in Section E.3.
Comparison (Cpr). We create 200 T/F pairs (400 queries total) focusing on comparing attributes or actions related to one or more human individuals across multiple images in a sequence. The ground truth is balanced at 100 True vs. 100 False.
Counting (Cnt). This subtask involves 120 nu-

merical questions asking for the number of people present or the frequency of certain actions in a sequence. Typical numeric answers range from 1 to 4, given the scope of each visual sequence.

Grouping (Grp). We provide 100 MC questions based on sequences containing at least three images, with at least two images featuring the same main "meta-human." The goal is to identify correct groupings of persons based on appearance, role, or action. As with *OC-Grp*, each question includes a "None" option as either the correct or a distractor choice.

Free-form (VID). We introduce 200 open-ended queries that focus on various person-centric aspects, such as identifying roles or describing activities. These questions allow more flexibility in model responses and assess the ability to generate context-relevant answers.

C.3 Annotation Quality and Agreement

As noted in the main text, three annotators reviewed all 3,060 question-answer pairs. An inter-annotator agreement study showed a high consensus rate of 98.74%, ensuring that the data is both accurate and consistent.

C.4 Summary

Our construction methodology ensures a balanced coverage of both object-centric and person-centric reasoning, as well as basic general cues such as element matching and tracking. The inclusion of multiple question types (T/F, MC, numerical, and open-ended) further promotes comprehensive eval-1020 uation of vision-language models. Figure 4 in the 1021 main paper illustrates the distribution of these sub-1022 tasks and their question-format breakdown. We 1023 believe that the richness and diversity of VLM²-1024 Bench make it a robust platform for advancing 1025 multimodal research. 1026

1031

1032

1033

1034

1035

1036

1039

1040

1041

1042

1043

1044

1046

1048

1049 1050

1051

1052

1053

1054

1056

1028

D.1 Chance-level

D

Baselines

1029In this part, we explain the calculation of chance-1030level accuracy for all tasks in our benchmark.

GC-Mat, GC-Trk. The Matching (Mat) and Tracking (Trk) tasks in General Cue (GC) follow a True-False (TF) paired-question format, where each pair consists of a positive question and a negative question:

• **Positive Question**: Derive from the correct *element* or *change*. The ground truth (GT) answer is True (T).

• Negative Question: Derive from the distractor *element* or *change*. The ground truth (GT) answer is False (F).

Positive Question:

"Is the answer 'the salad' correct for the given question: 'What object that was present in the first image is no longer visible in the second?'"

GT Answer: True (T)

Negative Question:

"Is the answer 'the ciabatta roll' correct for the given question: 'What object that was present in the first image is no longer visible in the second?'"

GT Answer: False (F)

Table 4: Example of True-False paired-question format in GC tasks.

During the construction of these questions, we ensure that the queried content originates from either the correct answer or a distractor answer. These elements are designed to be **independent** and identically distributed. Since each question in the pair has an independent 50% chance of being answered correctly, the expected accuracy under random guessing would be $P(\text{correct answer}) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} = 25\%$.

OC-Cpr, PC-Cpr. The OC-Cpr and PC-Cpr tasks utilize a True-False (TF) paired-question format where both questions in a pair originate from the same correct answer but are framed in two different ways:

- Positive Question: A direct affirmative statement that correctly represents the ground truth.
 1057
 1058
 1059
- Negative Question: A negated version of the positive question, often by inserting "not" after the verb.

An example is shown in Table 5.

Positive Question:

"Given the images, the claim 'The pets in these images are the same pet.' is right." GT Answer: **True** (**T**)

Negative Question:

"Given the images, the claim 'The pets in these images are not the same pet.' is right." GT Answer: False (F)

Table 5: Example of TF paired-question format in OC-Cpr and PC-Cpr tasks.

This construction aims to eliminate **language bias** by ensuring that the model does not favor one phrasing over another. For a language model that is free from bias, these two questions are **logically equivalent**—answering one correctly implies answering the other correctly as well. Consequently, under random guessing, the expectation is $P(\text{correct answer}) = \frac{1}{2} = 50\%$.

OC-Cnt, PC-Cnt. The calculation formulas for the accuracy of the chance-level accuracy are the same as in Section 3.1.

Under a pure random guessing strategy, the predicted answer \hat{N}_i is uniformly sampled from the set $\{1, 2, ..., L\}$, where L is the number of images (i.e., the sequence length for that instance). For a fixed sequence length L, we can compute the expected normalized accuracy E(L) by averaging over all possible ground-truth and guess pairs:

$$E(L) = 1 - \frac{1}{L^2} \sum_{N=1}^{L} \sum_{\hat{N}=1}^{L} w(L) \cdot \epsilon(N, \hat{N})^{-\alpha},$$
 108

where

$$\epsilon(N, \hat{N}) = \frac{|\hat{N} - N|}{\max(N - 1, L - N)}$$
 1084

and the weight is defined as

$$w(L) = \frac{L_{\max}}{L},$$
 1080

1083

1085

1060

1061

1063

1064

1065

1066

1067

1069

1070

1071

1072

1073

1075

1076

1078

1079

1080

1081

A question pair example is shown in Table 4.

1087	with $L_{\rm max} = 4$ being the maximum sequence
1088	length in our dataset.
1089	OC-Cnt Task: The OC-Cnt task exhibits the
1090	following distribution:
1091	• Length 2: 80 sequences (22.2%)
1092	• Length 3: 120 sequences (33.3%)
1093	• Length 4: 160 sequences (44.4%)
1094	Thus, the overall chance level accuracy is ob-
1095	tained as the weighted average: $Acc_{\text{OC-Cnt}} =$
1096	$\frac{80 E(2) + 120 E(3) + 160 E(4)}{360} \approx 34.88\%.$
1097	PC-Cnt Task: For the PC-Cnt task, the se-
1098	quence distribution is:
1099	• Length 2: 30 sequences (25.0%)
1100	• Length 3: 25 sequences (20.8%)
1101	• Length 4: 65 sequences (54.2%)
1102	Accordingly, the overall chance level accuracy is
1103	given by: $Acc_{PC-Cnt} = \frac{30 E(2) + 25 E(3) + 65 E(4)}{120} \approx$
1104	34.87%.
1105	D.2 Human-level
1106	To facilitate human participants in providing re-
1107	sponses to our questions, we integrated all model-
1108	prompted questions and answer choices into a
1109	graphical user interface (GUI), as illustrated in Fig-
1110	ure 6. This interface enabled participants to select

their answers conveniently, ensuring consistency

in data collection. We then gathered all responses

and conducted statistical analysis on the collected

human evaluations.

1111

1112

1113

1114



Figure 6: GUI for human-level testing

1138

1139

1140

1141

1142

E More details on Benchmark Construction

E.1 GC (General Cue) 1117

Manual Screening and Refine. Figure 7 demon-1118 strates the Graphic User Interface (GUI) we build 1119 for manually screening image editing data. 1120

The Pseudocode in Figure 8 Salient Sampling. 1121 and Table 6 displays the calculation process for the 1122 salient sampling score mentioned in Section 2.1. 1123

1124 Prompts for Pair-wise Answer Generation. Table 7 and 8 provides the complete prompts used 1125 to generate pair-wise answers for our evaluation 1126 tasks. The prompts were designed to instruct the 1127 language model to produce two distinct answers-a 1128 1129 positive (T) answer and a negative (F) answer—for each task. The dual-answer format is intended to 1130 capture both the expected response and its direct 1131 opposite, thereby offering a more balanced insight 1132 into the model's understanding. 1133

Question Templates. Table 9 and 10 list de-1134 tailed standard question templates for General Cue -1135 Matching and Tracking tasks, including the format 1136 instruction prompt. 1137

E.2 OC (Object-centric Cue)

Data Collection. To construct the dataset, we follow a structured approach to collect object-centric images, as illustrated in Figure 9. In total, we manually collected 320 images for objects.

Main Meta-Object Selection. We predefine 8 1143 types of common objects, with each type contain-1144 ing 5 meta-objects. For each meta-object, we col-1145 lect four images that represent the same object from 1146 different angles and scene conditions. 1147

Distractor Meta-Object Selection. To build 1148 meaningful object image sequences, we introduce 1149 visually distractive elements for each main meta-1150 object, referred to as "distractor meta-objects". 1151 Specifically, for each main meta-object, we col-1152 lect four additional images that belong to different 1153 but visually similar meta-objects within the same 1154 object category. These images are selected fol-1155 1156 lowing predefined visual cue confusion principles, ensuring that they provide meaningful challenges 1157 for vision language models. We ensure that each 1158 distractor image belongs to a different distractor 1159 meta-object, fundamentally guaranteeing that the 1160

count of different meta-objects in the final con-1161 structed sequence strictly follows our design. The 1162 principle of selecting distractor meta-objects is il-1163 lustrated in the outer ring of Figure 9. 1164

Image Sources. The images are gathered from 1165 various sources based on the nature of the objects: 1166

• Plush Objects: Images of plush toys are entirely sourced from the Jellycat website and its 1168 review sections, where diverse user-uploaded images provide a wide variety of object angles 1170 and scenes. 1171

1167

1169

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1202

1204

- Pet Objects: For the pet category of metaobjects, we source images from a combination of social media accounts of popular pet influencers' pet photography. We also include images of a ragdoll cat owned by one of the authors. As a result, this approach guarantees that each pet meta-object within the dataset belongs to the same individual cat or dog, minimizing variability unrelated to visual cue confusion.
- Other Objects: Most images are collected from Amazon product listings and review sections containing user-uploaded photos. A smaller portion of the dataset is curated using Google Lens image search, where specific visual distractive cues are used to retrieve and manually select images. The detailed visual cue principles guiding this selection process can be found in Figure 9.

Images Sequence Construction. The construction of image sequences in OC (a total of 360 sequences) follows the structure in Table 11. More specific details are listed below:

Two-Image Sequences (image_seq_len = 2)

- 1. Main Meta-Object Only (AA): Two images 1196 are randomly sampled from the same main 1197 meta-object. 40 sequences are constructed 1198 (one for each main meta-object). 1199
- 2. Main Meta-Object + Distractor Meta-1200 Object (AB): One image is randomly selected 1201 from the main meta-object, and one from the corresponding distractor meta-object. 40 se-1203 quences are constructed.
- Three-Image Sequences (image_seq_len = 3) 1205

	Image Filter App
	<image/>
	ID: 182759_3
	Progress: Item 68/227 (Remaining: 160)
	Edit Prompt:
	Keep (1) Skip (0) Save Now
Edit Prompt:	swap the dog's closed mouth with a yawning mouth
Select Task:	swap 🔽
Saliency of change (1-5):	3
Location of change (1-5):	3
Background complexity (1-5):	3
Time of spotting (1-5):	2
From coco (T/F):	
	Submit Scores

Figure 7: GUI.

Supposed you are looking at two images: Image 1: <<u>Cap_src></u> Image 2: <<u>Cap_edit></u> From Image 1 to Image 2, the change can be summarized as: <<u>P></u>

Table 6: Template for salient-score calculation.

#Task Description

Given the change between the first image and the second image, you need to generate four choices to the question "What new element can be observed in the second image that was not present in the first?" (this question varies based on the editing task). Remember, the choices' lengths should be similar. Additionally, your response should start with "Choices".

#Pair Design

In these two choices, you need to contain *only* the names of objects, but be specific:

1. Correct Answer (You need to infer the *only* from the Editing Information)

2. Distractor (You need to pick a random object *only* in the 'Description', but differ from the correct answer object)

#In-context example

Editing Information:

Add a katana held in the figure's left hand, angled downwards.

Description:

The image depicts a person dressed in traditional Japanese armor, standing in a misty, snowy landscape. The armor is detailed and appears to be made of metal, with various straps and buckles. The person is wearing a black mask that covers their entire face, adding to the mysterious and stealthy appearance. The background features stone lanterns and other traditional Japanese structures, which are partially obscured by the mist. The overall atmosphere is serene yet somewhat eerie, with the mist adding a sense of mystery and isolation. The scene suggests a historical or fantasy setting, possibly a samurai or ninja in a snowy, misty environment.

Choices:

Correct Answer: katana held Distractor: black mask **#Task <Original Edit Prompt> <VLM's Description>**

Table 7: Pair-aware answer generation prompt for OC-Mat

Task Description

Given the change between the first image and the second image, you need to generate four choices to the question "What key visual difference can be observed from the first image to the second image?". Remember, the choices' lengths should be similar. Additionally, your response should start with "Choices: " and must contain Correct Answer and Direct Reverse Answer.

Pair Design

In the two choices, you need to contain:

1. Correct Answer (You need to infer from the Editing Information)

2. Direct Reverse Answer (You need to infer from the Editing Information and change it to the opposite)

In-context example

Editing Information:

Swap the black ninja gloves with clean white gloves appropriate for serving.

Description:

The image depicts a person dressed in formal attire, standing in a doorway. The individual is wearing a black tuxedo with a white dress shirt and a black bow tie. They are holding a tray with several items on it. The tray contains a small glass container, a bottle, and a small white object, possibly a salt shaker or a similar item. The person is also wearing black gloves, which are typical for serving or formal dining scenarios. The background shows a wooden door with a brass hinge and a light-colored wall. The setting appears to be indoors, possibly in a house or a formal establishment.

Choices:

Correct Answer: The black ninja gloves were replaced with clean white gloves.

Direct Reverse Answer: The clean white gloves were replaced with black ninja gloves. **#Task**

<Original Edit Prompt> <VLM's Description>

Table 8: Pair-aware answer generation prompt for OC-Trk

Algorithm 1 Salient Score Computation

```
# cap_src: caption for the source image
     cap_edit: caption for the edited image
     T: template for constructing a paragraph
  # P: editing prompt
  input_text = concat(cap_src, cap_edit, T)
   in_tokens = tokenizer.encode(input_text)
   out_tokens = tokenizer.encode(P)
   log_sum = 0
   tokens = in_tokens
   # Model Forward Pass
11
   for i in range(1, len(out_tokens)):
       outputs = model(tokens)
logits = outputs.logits
14
15
       # Extract log probability of next token
probs = log_softmax(logits[0, -1, :])
16
17
       prob = probs[out_tokens[i]]
18
19
       log_sum += prob
20
         Update Input Sequence
21
       tokens = concat(tokens, out_tokens[i])
24 # Normalize the total log probability as the
  salient_score
salient_score = log_sum / len(out_tokens)
25
27 # Return: salient_score
```

Figure 8: Pseudocode for salient score computation.

1. Main Meta-Object Only (AAA): Three images are randomly sampled from the same main meta-object. 40 sequences are constructed.

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225 1226

1227

1228

1229

- 2. Main Meta-Object + Distractor Meta-Object (AAB): Two images are selected from the main meta-object, and one from the distractor meta-object. The order of images is shuffled. 40 sequences are constructed.
- 3. Main Meta-Object + Distractor Meta-Objects (ABC): One image is selected from the main meta-object, while two are selected from different distractor meta-objects. 40 sequences are constructed.
- Four-Image Sequences (image_seq_len = 4)
 - 1. Main Meta-Object Only (AAAA): All four images are sampled from the same main meta-object and shuffled. 40 sequences are constructed.
- 2. Main Meta-Object + Distractor Meta-Object (AAAB): Three images are sampled from the same main meta-object, while one is selected from a distractor meta-object. 40 sequences are constructed.
- 12303. Main Meta-Object + Distractor Meta-1231Objects (AABC): Two images are selected

GC-Mat Positive Question:

"Is the answer 'correct element' correct for the given question: 'What new element can be observed in the second image that was not present in the first?'" GT Answer: **True** (**T**)

GC-Mat Negative Question:

"Is the answer 'distractor element' correct for the given question: 'What new element can be observed in the second image that was not present in the first?'" GT Answer: False (F)

Table 9: GC-Mat True-False paired-question

GC-Trk Positive Question:

"Is the answer 'correct change' correct for the given question: 'What key visual change can be observed from the first image to the second image?'"

GT Answer: True (T)

GC-Trk Negative Question:

"Is the answer 'distractor change (reversed process)' correct for the given question: 'What key visual change can be observed from the first image to the second image?'" GT Answer: False (F)

Table 10: GC-Trk True-False paired-question

from the main meta-object, while two are se-
lected from different distractor meta-objects.1232
1233
123440 sequences are constructed.1234

4. Main Meta-Object + Distractor Meta-
Objects (ABCD): One image is selected from
the main meta-object, while three are selected
from different distractor meta-objects. 40 se-
quences are constructed.1235
1236
1238

Question Templates.Table 12, 13 and 14 list1240detailed standard question templates (with format1241instructions) for the Object-centric Cue task, in-1242cluding 3 subtasks: Comparison (cpr), Counting1243(Cnt), and Grouping (Grp).1244



Figure 9: The overview of the structured design of the Object-centric Cue (OC) images. **Central Layer (Main Meta-Objects)**: The innermost circle represents the predefined **8 object categories**, which serve as the foundation for our dataset. These categories include *Pet, Plush, Bag, Book, Cup, Shirt, Shoes, and Toy*. Each category consists of 4 main meta-objects. **Middle Layer (Example Meta-Objects within Each Category)**: Each segment surrounding the center showcases a representative **main meta-object** within its category. These meta-objects serve as core instances for data collection. For example, the *Pet* category includes *Cat* and *Dog*, while the *Bag* category includes *Backpack, Schoolbag* and *Fashion Bag*. **Outer Layer (Distractor Meta-Objects & Visual Cue Distraction Principles)**: The outermost ring presents 1 out of 4 **distractor meta-objects** specifically selected to create challenging image sequences. Each distractor meta-object shares one or more **distractive visual cues** with its corresponding main meta-object.

Num	Src	Process of Image Sequences Construction	Cpr	cnt	Grp
2	AA	2 images from the same object O_i , randomly sampled as $\mathcal{I}_{O_i} = \{I_i, I_j\}$, and shuffled.	Т	2	-
2	AB	1 image I_i from \mathcal{I}_{O_i} and 1 image $I_{\neg i}$ from distractor set $\mathcal{I}_{\neg O_i}$, randomly shuffled.	F	1	-
3	AAA	3 images from the same object O_i , randomly sampled as $\mathcal{I}_{O_i} = \{I_i, I_j, I_k\}$, and shuffled.	Т	3	-
3	AAB	2 images from the same object O_i , randomly sampled as $\mathcal{I}_{O_i} = \{I_i, I_j\}$ and 1 $I_{\neg i}$ from distractor set $\mathcal{I}_{\neg O_i}$, randomly shuffled.	F	2	$[I_i, I_j]$
3	ABC	1 images from the same object O_i , randomly sampled as $\mathcal{I}_{O_i} = \{I_i\}$ and 2 images $\{I_{\neg i}, I_{\neg j}\}$ from distractor set $\mathcal{I}_{\neg O_i}$, randomly shuffled.	F	3	[]
4	AAAA	4 images from the same object O_i , randomly sampled as $\mathcal{I}_{O_i} = \{I_i, I_j, I_k, I_p\}$, and shuffled.	Т	4	-
4	AAAB	3 images from the same object O_i , randomly sampled as $\mathcal{I}_{O_i} = \{I_i, I_j, I_k\}$ and 1 image $I_{\neg i}$ from distractor set $\mathcal{I}_{\neg O_i}$, randomly shuffled.	F	2	$[I_i, I_j, I_k]$
4	AABC	2 images from the same object O_i , randomly sampled as $\mathcal{I}_{O_i} = \{I_i, I_j\}$ and 2 images $\{I_{\neg i}, I_{\neg j}\}$ from distractor set $\mathcal{I}_{\neg O_i}$, randomly shuffled.	F	3	$[I_i, I_j]$
4	ABCD	1 images from the same object O_i , randomly sampled as I_i and 3 images $\{I_{\neg i}, I_{\neg j}, I_{\neg k}\}$ from distractor set $\mathcal{I}_{\neg O_i}$, randomly shuffled.	F	3	[]

Table 11: Summary of Multi-Images Sequence Construction for Object-centric Cue (OC) Tasks

OC-Cpr Positive Question:

"Answer the following question according to this rule: You only need to provide *ONE* correct answer with 'T' (True) or 'F' (False). Only reply with 'T' or 'F'. The Question is: Given the images, the claim 'The {obj}s in these images are the same {obj}.' is right."

GT Answer: **True** (**T**)

OC-Cpr Negative Question:

"Answer the following question according to this rule: You only need to provide *ONE* correct answer with 'T' (True) or 'F' (False). Only reply with 'T' or 'F'. The Question is: Given the images, the claim 'The {obj}s in these images are not the same {obj}.' is right."

GT Answer: False (F)

Table 12: OC-Cpr True-False paired-question

OC-Cnt Question:

"Answer the following question according to this rule: You only need to provide *ONE* correct numerical answer. For example, if you think the answer is '1', your response should only be '1'. The Question is: How many different {obj}s are there in the input images?" GT Answer: "N" (e.g., "1", "2", etc.)

Table 13: OC-Cnt Numerical Counting Question

OC-Grp Question:

"Answer the following question according to this rule: You only need to provide *ONE* correct answer with the corresponding letter. For example, if you think the correct answer is 'B) 1 and 2', your response should only be 'B) 1 and 2'. The Question is: Which images show the same {obj} in the input images? Choices: A) 1 and 3; B) None; C) 2 and 3; D) 1 and 2." GT Answer: "A) 1 and 3" (Example Answer)

Table 14: OC-Grp Multiple-Choice Grouping Question

E.3 PC (Person-centric Cue)

1245

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258 1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1270

Data Collection. We collect images of *meta-humans* mainly from https://www.imdb.com/ and some are from the actor or actress's social media.

Main Meta-human Selection. Our dataset is evenly distributed across different racial groups (Asian, Black, and White) and genders (Male and Female). For every race-gender combination, we select five main meta-humans, each contributing four images, yielding a total of 120 images.

To ensure consistency, all selected individuals are within a similar age range, preventing significant age-related facial changes that could interfere with identity recognition. Additionally, each actor's appearance remains relatively consistent in terms of makeup and overall styling, ensuring that different images of the same meta-human retain distinct yet comparable visual cues (e.g. face shape, eye spacing, nose structure, and lip contours). By preserving these features, we avoid manipulating a single individual's visual cues that could potentially mislead VLMs. Rather, we ensure that the evaluation genuinely tests whether the model can visually link matching cues to recognize the same or different individuals without prior identity knowledge.

Distractor Meta-human Selection. To intro-1271 duce challenging distractors in our sequences, we 1272 compute the CLIP embedding for every image and 1274 store these embeddings in a reference base. When a distractor image is needed, we perform an image-1275 to-image similarity search within this base to iden-1276 tify the most visually similar image that originates 1277 from a different meta-human. This fine-grained 1278

matching ensures that the distractor image closely1279resembles the main meta-human's image, leading1280to more challenging image sequences.1281

Discussion on Why Objects Require Dedicated 1282 Distractors, While Humans Do Not. In object-1283 centric tasks, objects are categorized into eight 1284 distinct types, with substantial differences among 1285 different types (e.g. pets and bags). Therefore, 1286 each main meta-object requires dedicated distrac-1287 tors from the same object type to ensure meaningful 1288 comparisons. In contrast, humans belong to a sin-1289 gle category, meaning that any meta-human can 1290 serve as a distractor for another. Given that we 1291 compute CLIP embeddings to select visually sim-1292 ilar distractors, the constructed image sequences 1293 already present a significant challenge without the 1294 need for type-specific distractors. We also ensure 1295 diversity by selecting five main meta-humans for 1296 each race-gender pair, providing a sufficiently large 1297 pool from which to choose suitable distractors. Cor-1298 responding to our hypothesis, in the final curated se-1299 quences, most distractor meta-humans chosen were 1300 of the same race or gender as the main meta-human. 1301 Additionally, as shown in Table 1, these curated im-1302 age sequences along with our designed questions 1303 effectively challenge tested models, revealing their 1304 limited performances in visually linking matching 1305 cues on person-centric data. 1306

Images Sequence Construction. The construction of image sequences in PC (a total of 260 sequences) follows the structure in Table 15. More specific details are listed below: 1307

1308

1309

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

Two-Image Sequences (image_seq_len = 2)

- 1. Main Meta-Human Only (PP): Two images are randomly sampled from the same main meta-human. 50 sequences are constructed.
- 2. Main Meta-Human + Distractor Meta-Human (PQ): One image is randomly selected from the main meta-human, and one from a distractor meta-human. The order of images is shuffled. 50 sequences are constructed.

Three-Image Sequences (image_seq_len = 3)

1. Main Meta-Human Only (PPP): Three images are randomly sampled from the same1322main meta-human. 20 sequences are constructed.1324

2. Main Meta-Human + Distractor Meta-Human (PPQ): Two images are selected from the main meta-human, and one from a single distractor meta-human. The order of images is shuffled. 30 sequences are constructed.

1326

1327

1328

1329

1331

1332

1333

1336

1337

1338

1339

1340

1341

1342

1344

1345

1346

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

1363

1364

1365

1366

1368

1370

1371

 Main Meta-Human + Distractor Meta-Humans (PQR): One image is selected from the main meta-human, while the other two come from distinct distractor meta-humans. The order is shuffled. 10 sequences are constructed.

Four-Image Sequences (image_seq_len = 4)

- 1. Main Meta-Human Only (PPPP): All four images are sampled from the same main meta-human. The order is shuffled. 30 sequences are constructed.
- Main Meta-Human + Distractor Meta-Human (PPPQ): Three images are sampled from the main meta-human, while one is selected from a single distractor meta-human. The order is shuffled. 20 sequences are constructed.
- Main Meta-Human + Distractor Meta-Humans (PPQR): Two images are selected from the main meta-human, while two are selected from distinct distractor meta-humans. The order is shuffled. 20 sequences are constructed.
 - 4. Main Meta-Human + Distractor Meta-Humans (PQRS): One image is selected from the main meta-human, while three are selected from distinct distractor meta-humans. The order is shuffled. 30 sequences are constructed.

Video Construction. The video data for this benchmark is manually collected from Shutterstock². We selected ten common activity categories that an individual can perform: clean, cook, drink, exercise, listen, play, read, ride, walk, and work. For each category, we curated 10 sets of candidate video pairs, and each set consists of two videos.

To ensure motion consistency and length diversity, we carefully structured the final videos by concatenating clips while keeping the total duration within the **0-100**s time range. Figure 10 displays the sketch of concatenated video length distribution. The final compositions followed two formats:

- *P*->¬*P* format: A direct concatenation of two distinct clips (same length for each clip).
- P->¬P->P format: A sequence where the first clip and the third clip are sampled from the same candidate video, while the second clip is sampled from the second candidate video. (same length for the three clips) 1378



Figure 10: Distribution of video duration.

Regardless of the different default sampling1379methods for our baseline models in Table 16, both1380 $P \rightarrow \neg P$ and $P \rightarrow \neg P \rightarrow P$ formats ensure that every video clip has frames included in the sampling1382process:1383

- Uniform Sampling (8/16 frame): Each clip
 contributes a proportionate number of frames
 based on the total video length. Since in one
 concatenated video, all the sampled clips are
 the same length, this method guarantees at
 least 2 frames for each clip can be sampled as
 model input frames.
- **FPS Sampling (1fps)**: Since frames are sampled at a fixed rate, the structure of $P \rightarrow \neg P$ and $P \rightarrow \neg P \rightarrow P$ ensures that each clip is present long enough for multiple frames to be captured, regardless of its placement in the sequence. 1396

²https://www.shutterstock.com

Num	Src	Process of Image Sequences Construction	Cpr	cnt	Grp
2	PP	2 images from the same person P_i , randomly sampled as $\mathcal{I}_{P_i} = \{I_i, I_j\}$, and shuffled.	Т	2	-
2	PQ	1 image I_i from \mathcal{I}_{P_i} and 1 image $I_{\neg i}$ from distractor set $\mathcal{I}_{\neg P_i}$, randomly shuffled.	F	1	-
3	PPP	3 images from the same person P_i , randomly sampled as $\mathcal{I}_{P_i} = \{I_i, I_j, I_k\}$, and shuffled.	Т	3	-
3	PPQ	2 images from the same person P_i , randomly sampled as $\mathcal{I}_{P_i} = \{I_i, I_j\}$ and 1 $I_{\neg i}$ from distractor set $\mathcal{I}_{\neg P_i}$, randomly shuffled.	F	2	$[I_i, I_j]$
3	PQR	1 image from the same person P_i , randomly sampled as $\mathcal{I}_{P_i} = \{I_i\}$ and 2 images $\{I_{\neg i}, I_{\neg j}\}$ from distractor set $\mathcal{I}_{\neg P_i}$, randomly shuffled.	F	3	0
4	PPPP	4 images from the same person P_i , randomly sampled as $\mathcal{I}_{P_i} = \{I_i, I_j, I_k, I_p\}$, and shuffled.	Т	4	-
4	PPPQ	3 images from the same person P_i , randomly sampled as $\mathcal{I}_{P_i} = \{I_i, I_j, I_k\}$ and 1 image $I_{\neg i}$ from distractor set $\mathcal{I}_{\neg P_i}$, randomly shuffled.	F	2	$[I_i, I_j, I_k]$
4	PQQR	2 images from the same person P_i , randomly sampled as $\mathcal{I}_{P_i} = \{I_i, I_j\}$ and 2 images $\{I_{\neg i}, I_{\neg j}\}$ from distractor set $\mathcal{I}_{\neg P_i}$, randomly shuffled.	F	3	$[I_i, I_j]$
4	PQRV	1 image from the same person P_i , randomly sampled as I_i and 3 images $\{I_{\neg i}, I_{\neg j}, I_{\neg k}\}$ from distractor set $\mathcal{I}_{\neg P_i}$, randomly shuffled.	F	3	[]

Table 15: Summary of Multi-Images Sequence Construction for Person-centric Cue (PC) Tasks

Model Name	Uniform (8/16)	FPS (1fps)
LLaVA-OneVision-7B	 Image: A second s	×
LLaVA-Video-7B	✓	×
LongVA-7B	✓	×
mPLUG-Owl3-7B	✓	×
Qwen2-VL-7B	×	1
Qwen2.5-VL-7B	×	1
InternVL2.5-8B	✓	×
InternVL2.5-26B	\checkmark	×
GPT-4o	 Image: A second s	×

Table 16: Comparison of Different Sampling Methods

1397Thus, by maintaining the integrity of each clip's1398temporal structure, both $P \rightarrow \neg P$ and $P \rightarrow \neg P \rightarrow P$ 1399formats effectively ensure that every clip con-1400tributes frames to the final sampled frame input1401for all models.

1402

1403

1404

1405

1406

1407

Question Templates. Table 17, Table 18, Table 19, and Table 20 present the detailed standard question templates for the Person-centric Cue task, covering the four subtasks: Comparison (PC-Cpr), Counting (PC-Cnt), Grouping (PC-Grp), and Video Identity Description (PC-VID).

PC-Cpr Positive Question:

"Answer the following question according to this rule: You only need to provide *ONE* correct answer with 'T' (True) or 'F' (False). Only reply with 'T' or 'F'. The Question is: The individuals in these images are the same person." GT Answer: **True (T)**

PC-Cpr Negative Question:

"Answer the following question according to this rule: You only need to provide *ONE* correct answer with 'T' (True) or 'F' (False). Only reply with 'T' or 'F'. The Question is: The individuals in these images are not the same person." GT Answer: False (F)

Table 17: PC-Cpr True-False paired-question

PC-Cnt Question:

"Answer the following question according to this rule: You only need to provide *ONE* correct numerical answer. For example, if you think the answer is '1', your response should only be '1'. The Question is: How many distinct individuals are in the input images?" GT Answer: 2 (Example Answer)

Table 18: PC-Cnt Numerical Counting Question

PC-Grp Question:

"Answer the following question according to this rule: You only need to provide *ONE* correct answer with the corresponding letter. For example, if you think the correct answer is 'B) 2 and 3', your response should only be 'B) 2 and 3'. The Question is: Which images correspond to the same person in the input images? Choices: A) None; B) 2 and 3; C) 1 and 3; D) 1 and 2." GT Answer: D) 1 and 2 (Example Answer)

Table 19: PC-Grp Multiple-Choice Grouping Question

PC-VID Question:

"Give a comprehensive description of the whole video, prioritizing details about the individuals in the video."

Table 20: PC-VID Video Describing Question

1409

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429 1430

1431

1432

1433

1434

1435

1436 1437

1438

1439

1440

1441

1442

1443

1444

1445 1446

F More details on Prompting Approaches

F.1 Prompts for LLM-as-Evaluator

When models answer our free-form PC-VID questions, their responses are evaluated by GPT-40 using the scoring prompts detailed in Tables 21 and 22. Specifically, for videos following a $\mathcal{P} \rightarrow \neg \mathcal{P}$ sequence, GPT-40 assesses whether the model explicitly distinguishes that the first individual (\mathcal{P}) and the second individual ($\neg \mathcal{P}$) are different. In this case, if the model successfully makes this distinction, it receives a score of 1; otherwise, it is given a score of 0.

For videos that exhibit a $\mathcal{P} \to \neg \mathcal{P} \to \mathcal{P}$ (PQP) pattern, the evaluation is more nuanced. The evaluator model (GPT-40) checks two aspects: (1) whether the model correctly identifies that there are two distinct individuals (i.e., \mathcal{P} and $\neg \mathcal{P}$), and (2) whether the model explicitly recognizes that the final appearance belongs to the same individual as the first (\mathcal{P}). A perfect identification of both aspects yields a score of 2, while correctly distinguishing the individuals without explicitly linking the final appearance to the first results in a score of 1. If the model fails to distinguish between the individuals, a score of 0 is assigned.

F.2 Prompting Approaches for Probing on VLM²-Bench

CoT (CoT-normal). The normal version of the Chain-of-Thought prompt is shown in Table 23. We simply require the model to think 'step-by-step' to ensure self-reflection and self-correction, as well as the transparent thinking process.

CoT-special for GC. Table 24 shows a special version of the Chain-of-Thought prompt. According to the task features, we carefully analyze how a human being approaches and visually links matching cues for questions in GC, then curate this prompt as an imitation of the human visual linking process.

VP-grid for GC. Figure 11 displays a complete 1447 version of Visual Prompting with Grid assistance 1448 (VP-grid). Here we follow (Lei et al., 2024) to 1449 print a set of dot matrix onto the input image, ac-1450 1451 companied by the image order dimension concatenated with Cartesian coordinates as (image order 1452 index, colum index), row index). In the detailed tex-1453 tual prompt design, we also integrated references 1454 and explanations for the grids, allowing VLMs to 1455

leverage this visual assistance as spatial and visual matching references.

1456

1457

1458

1459

1460

1461

1462

1463

VP-zoom-o for OC. In Figure 12, we demonstrate the visual prompting process for OC. We leverage the Grounded-SAM (Ren et al., 2024) model to detect bounding boxes for objects based on their types then crop the "zoomed-in" objects as the image input for further VQA pairs.

VP-zoom-p for PC.The visual prompting pro-
cess in similar to that of OC (Figure 13). We use a
face detection model (Geitgey, 2016) to "zoom in"1465
1465on the individual's face and occlude other irrelevant
information.1466

#Task

You are evaluating a model's ability to accurately distinguish between two different individuals, P and Q, who appear sequentially in a video (first P, then Q). Given a description, your task is to determine if the model explicitly identifies that the first person (P) and the second person (Q) are different individuals.

#Return Format

You only need return a number after "Score:". If you think the model correctly identifies that the two appearances belong to different individuals, return "Score: 1". If you think the model fails to explicitly state that there are two different individuals, return "Score: 0".

#Description

<Model's Description>

Table 21: Scoring prompt for *VID* (when video belongs to category of $P \rightarrow \neg P$).

#Task

You are evaluating a model's ability to accurately distinguish between two different individuals, P and Q, who appear sequentially in a video following an PQP pattern (first P, then Q, then P again). Given a description, your task is to determine whether the model explicitly identifies that: (1) P and Q are different individuals, and (2) The person in the final scene is the same as the first (P). **#Return Format**

You only need return a number after "Score:".

(1) If the model correctly describes that the video follows an PQP sequence, explicitly recognizing that the first and last appearances belong to the same person (P), while the middle appearance is a different person (Q), return "Score: 2".

(2) If the model correctly identifies that there are two different people in the video (P and Q) but does not explicitly mention that the last scene returns to P, return "Score: 1".

(3) If the model fails to recognize that two different individuals appear (e.g., treats all appearances as the same person or does not distinguish between P and Q), return "Score: 0".

#Description

<Model's Description>

Table 22: Scoring prompt for *VID* (when video belongs to category of $P \rightarrow \neg P \rightarrow P$).

<Question>

Let's think 'step by step' to answer this question, you need to output the thinking process of how you get the answer.

Table 23: CoT prompt for GC (here we denote as CoT-normal to distinguish it from the CoT-special in Table 24 that specifically designed for GC), OC, and PC.

<Question>

Use the following 4 steps to answer the question:

Step 1. Understand the Question

- Identify the question's purpose.
- Check for any format requirements.

Step 2. Perceive (List Elements)

- List every details in each image respectively.
- Note positions and attributes of elements.

Step 3. Connect (Compare & Reason)

- Compare corresponding elements in each image.
- List all the unchanged elements and the changed element.

Step 4. Conclude (Answer the Question)

Table 24: CoT-special specifically designed for GC.



<Question>

Here's the instruction you need to strictly follow to approach this question:

Two images are provided, each overlaid with a grid of dots arranged in a matrix with dimensions h by w. Each dot on this grid is assigned a unique set of three-dimensional coordinates labeled as (t, x, y). The first coordinate, "t," distinguishes the two images— "1" for the first image, "2" for the second. The remaining coordinates, "x" and "y," specify each dot's location, where within any column x increases from top to bottom, and within any row y increases from left to right.

This labeling system is intended to help you identify, reference, connect, and compare objects across both images. Now, use the following 4 steps to answer the question.

Step 1. Understand the Question - Identify the question' s purpose. - Check for any format requirements.

Step 2. Perceive (List Elements and coordinates) - For all the objects in the 'Options' of the question, identify them in each image separately, double check their existence. If the object exists then output its nearest coordinates. - Output format like 'Image1: apple at coordinates (1, 2, 3)... Image2: banana at coordinates (2, 4, 5)'

Step 3. Connect (Compare & Reason) - Use the grid coordinates to connect objects across the two images, observing any similarities or differences at the same (x, y) positions.

Step 4. Conclude (Answer the Question) - If a specific output format is required (e.g., "MY_ANSWER: ..."), follow it exactly. Include the transparent thinking process in your answer, and make sure you output the final *ONE* answer after 'MY_ANSWER:', just like 'MY_ANSWER: D) '

Figure 11: An illustration of how VP-grid works for GC.



Figure 12: An illustration of how VP-zoom-o works for OC.



Figure 13: An illustration of how VP-zoom-p works for PC.

G Case Study

1469

1470

1471

1472

1473

1474

1475

This section focuses on how various prompting techniques influence model performance, highlighting their successes and limitations across different models.

G.1 Case for CoT-special prompting in General Cue (GC) Task

We observe that the CoT-special prompt boosts 1476 InternVL2.5-8B's performance by over 25% than 1477 the standard query in both Matching and Track-1478 ing tasks for General Cue. While for the traditional 1479 CoT-normal prompting technique, this boost is only 1480 13%. The CoT-special prompt (Table 24) directs 1481 the model through four explicit steps: understand-1482 ing the question, perceiving (listing elements), con-1483 1484 necting (comparing and reasoning), and concluding. This structured approach mirrors the human 1485 process of visual matching and is effective even for 1486 a rather smaller model like InternVL2.5-8B, which 1487 might otherwise struggle with the ambiguity of a 1488

complex generic step-by-step instruction (which we will discuss later in the next Subsection G.2). 1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

For example, in the provided InternVL2.5-8B response Figure 14, the model correctly executes the following: In Step 2, it identifies critical details such as "Vase with flowers on the table" and "Chandelier above" in Image 1, while noting the absence of the vase in Image 2. In Step 3, it systematically compares the two images, highlighting that while many elements remain unchanged (e.g., the chandelier, kitchen area, bowl of fruit, window), the removal of the vase is the key difference. Finally, in Step 4, the model concludes that the statement "The vase on top of the table was removed" accurately describes the visual change, thereby arriving at the correct answer.

This detailed, multi-step breakdown not only1505ensures that all pertinent visual cues are captured1506and processed but also reduces errors by structur-1507ing the logical flow of reasoning. The CoT-special1508prompt's explicit instructions help InternVL2.5-8B1509

1510align visual information with textual descriptions1511more effectively, thus enhancing overall perfor-1512mance. Compared to the less specific CoT-normal1513prompt—which may leave the model with gaps in1514reasoning—the CoT-special prompt provides clear,1515task-specific guidance that is essential for complex1516visual reasoning tasks, as evidenced by the substan-1517tial performance improvement.

1518 G.2 Case for VP-grid in General Cue Task

1519

1520

1521

1523

1524

1525

1527

1528

1530

1531

The VP-grid (Visual Prompting with Grid assistance) method enhances visual matching in General Cue tasks by overlaying a dot matrix grid onto the input image. Each dot is annotated with a three-dimensional coordinate tuple, (*image order index, column index, row index*), where the first dimension distinguishes the sequence of images (e.g., the first image is indexed as 1 and the second as 2). This grid is further supported by detailed textual descriptions that clarify the coordinate system, enabling Vision-Language Models (VLMs) to use these cues for spatial and visual matching.

A example failure case in VP-grid. However, 1532 this approach does not yield consistent improve-1533 ments across all models. For instance, the 1534 Qwen2.5-VL-7B model demonstrates a significant 1535 performance drop—nearly 20%—when using VP-1536 grid. An example failure case is in Figure 15. 1537 Our analysis reveals that although the model cor-1538 rectly identifies visual elements (e.g., a pedestrian 1539 1540 with a high-visibility vest at coordinates (2, 5, 3)), it fails to properly interpret the image sequence. 1541 Specifically, the model incorrectly associates the 1542 coordinates (2, 5, 3) with the first image, rather than the second, despite the explicit definition pro-1544 1545 vided in the textual prompt. This misinterpretation leads to erroneous linking of visual matching cues 1546 and subsequent faulty reasoning. We suspect that 1547 the underlying issue is the limited semantic com-1548 prehension capability of the relatively smaller 7B 1549 model, which struggles with complex, predefined 1550 spatial instructions and visual assistance. 1551

1552A example of success case in VP-grid. In con-1553trast to models that often misinterpret or neglect1554spatial cues provided by VP-grid—leading to errors1555such as mismatching image indices—GPT-40 suc-1556cessfully leverages these visual prompts to achieve1557correct visual-textual alignment. In the example at1558Figure 16, the model identifies the cat's nose at co-1559ordinates (1, 2, 4) in the first image and at (2, 2, 4)

in the second image, enabling it to accurately capture the change in the visual attribute (from a lighter pink to a darker black).

This success stems from several key aspects of GPT-4o's processing capabilities:

- 1. **Precise Disambiguation of Image Order:** The VP-grid explicitly encodes image order, which GPT-40 uses to differentiate between multiple images. This prevents the common error of conflating spatial information from distinct images—a problem seen in smaller models.
- 2. **Robust Visual Matching in space:** With clear coordinate annotations, the model effectively locates and compares the same physical regions across images. In this case, the exact correspondence between the cat's nose in different images is recognized, which is crucial for detecting subtle visual changes.
- 3. **Structured Reasoning Process:** GPT-4o adheres to a well-defined reasoning sequence in our textual guidance(perception, connection, and conclusion). By systematically linking the provided grid coordinates with the textual descriptions, it is able to deduce the key visual change accurately.

Implications on Model Scale. Our analysis suggests that the enhanced performance of GPT-40 with VP-grid can be attributed to its larger model capacity. Although the detailed architecture of GPT-40 is proprietary, its ability to process complex multi-modal prompts implies that:

- Enhanced Semantic Understanding: Larger models are inherently better at comprehending intricate, structured prompts that combine visual and textual information. This results in a more nuanced interpretation of spatial cues.
- Superior Visual-Textual Alignment: With greater capacity, GPT-40 can integrate and correlate the detailed spatial data (visual assistance) from the VP-grid with the corresponding textual descriptions, minimizing the risk of mis-association or errors.
- Effective Handling of Complexity: The advanced reasoning capabilities of larger models enable them to navigate the additional complexity introduced by VP-grid without suffering from the side effects seen in smaller models. This ensures that the additional spatial



Figure 14: Case study on why CoT-special leads to performance improvement.

1609

1610

1611

1612

1649 1650

1655

1656

1658

guidance improves performance rather than causing confusion.

The success of GPT-40 in utilizing the VP-grid approach demonstrates that model scale plays a critical role in effectively integrating complex visual and textual cues. By accurately disambiguating image order and performing precise spatial matching, GPT-40 not only avoids the pitfalls encountered by smaller models but also benefits significantly from the additional visual assistance, leading to an overall performance improvement of approximately 10%.

Case for CoT prompting in **G.3 Object-centric Cue Task**

The task design for Object-centric cue (OC) and person-centric cue (PC) requires multiple images (more than 2) as sequence input. We observe that, unlike General Cue (GC) tasks where models are required to link instance-level cues, OC tasks demand that models group similar objects based on fine-grained visual features. As illustrated in Figure 5b, models using the CoT approach sometimes struggle to provide a comprehensive overview of vision-based cues across a sequence of images.

A detailed case in Figure 17 is provided by InternVL2.5-26B's response. The ground truth and Vanilla responses correctly identify that there is no grouping for the same meta-object in the sequence, with the answer 'D) None'. In the C o T response, the model states: "The second and third images both have dinosaurs wearing sunglasses". Although the description here is true, its ambiguity and lack of detailed coverage lead the model to incorrectly select option C) 2 and 3, rather than the correct option D) None. Because if we take a closer look at the design on the backpack in image 3, the dinosaur with sunglasses is actually holding a keyboard instead of a skateboard in image 2. This is a distractive visual matching cue we intend to capture during the distractor meta-object selection. This major difference should have prevented models from grouping image 2 and image 3 together.

According to our findings, this misgrouping occurs for two main reasons:

1. Insufficient Overview of Visual Cues: The CoT prompt does not force the model to systematically verify all critical details across multiple images. As a result, the model overlooks nuanced differences, such as the design discrepancy on the backpack in image

3, where the dinosaur holds a keyboard rather than a skateboard.

1659

2. Variability in Descriptive Language: The 1661 open-ended language generated by the CoT 1662 approach can lead to inconsistent descriptions. 1663 In this case, the model generalized the visual 1664 cue of a "dinosaur design" without capturing 1665 the specific attribute (i.e., the object the di-1666 nosaur is holding), which is crucial for correct grouping.

Thus, the lack of structured guidance in the CoT prompt leads to the dropping or misinterpretation 1670 of critical cues, resulting in incorrect grouping deci-1671 sions for multi-image sequences in OC tasks. This 1672 analysis underscores the importance of more de-1673 tailed structured intermediate reasoning strategies, 1674 such as those provided by a tailored CoT-special 1675 prompt, to ensure that all relevant visual details are 1676 captured and compared accurately. 1677



Figure 15: Case study on why VP-grid leads to performance degradation.



Figure 16: Case study on why VP-grid leads to performance improvement for GPT-40



Figure 17: Case study on why CoT leads to performance degradation.