



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

Jianshu Zhang<sup>♥\*</sup>, Dongyu Yao<sup>♣\*</sup>, Renjie Pi<sup>♡</sup>, Paul Pu Liang<sup>♦</sup>, Yi R. (May) Fung<sup>♡</sup>

<sup>♡</sup>Hong Kong University of Science and Technology

<sup>♣</sup>Carnegie Mellon University

<sup>♦</sup>Massachusetts Institute of Technology

jianshu.zhang777@gmail.com rainy@cmu.edu rpi@ust.hk

ppliand@mit.edu yrfung@ust.hk

## 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 twelve VLMs, along with further analysis of various language-side 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. 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.<sup>1</sup>

## 1 Introduction

Humans constantly link matching visual cues to navigate and understand their environment. For instance, we can determine whether objects, and 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

<sup>\*</sup>These authors contribute to this work equally.

<sup>1</sup>Project page: <https://vlm2-bench.github.io/>.

<sup>♣</sup>Work was done while student was an intern at HKUST.

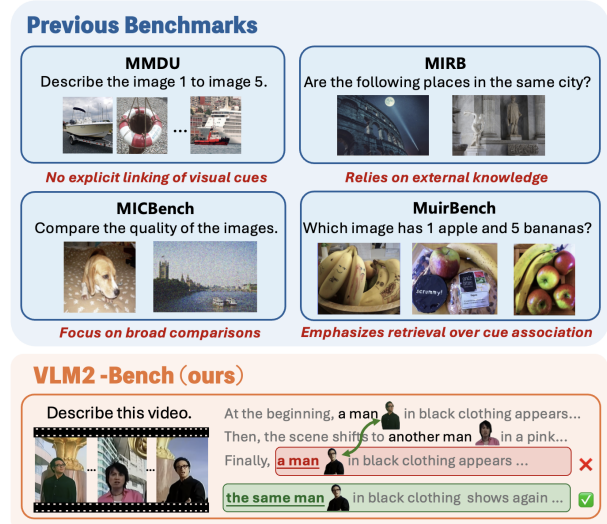


Figure 1: **Previous benchmarks** fail to assess the ability to link matching visual cues, whereas our **VLM2-Bench** explicitly tests this ability, as shown in the example where the model need to identify the reappearance of the same person by linking visual cues, like facial features or clothing, across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

us. However, while current vision-language models (VLMs) (Chen et al., 2024b; Li et al., 2024b; Zhang et al., 2024b; Team, 2025) have demonstrated extensive knowledge and expanded their capabilities from single-image understanding to handling multiple images and videos, *whether they 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

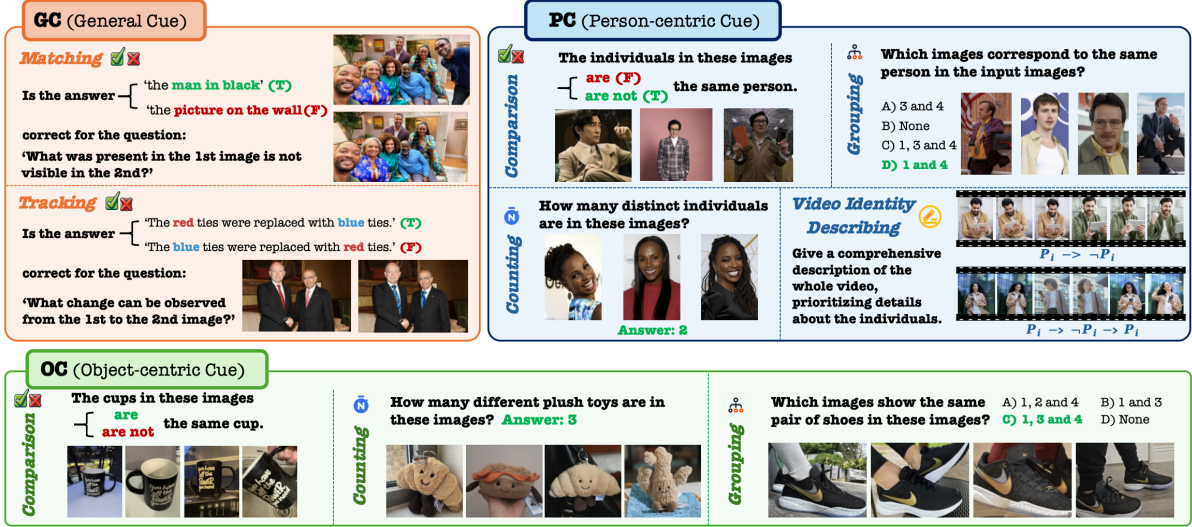


Figure 2: Overview of **VLM2-Bench**. The benchmark is categorized into three subsets based on visual cues: GC (General Cue), OC (Object-centric Cue), and PC (Person-centric Cue), each comprising multiple subtasks. To comprehensively evaluate VLMs’ ability to visually link matching cues, the benchmark includes diverse question formats—T/F ✓✗, multiple-choice ②, numerical ②, and open-ended ②—ensuring a comprehensive evaluation.

(d) focus on retrieval-based tasks rather than evaluating the direct association of visual cues across different visual contexts (Wang et al., 2024a).

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 tasks.

We conduct a comprehensive evaluation of 8 open-source models and 3 commercial models 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, even the most advanced commercial models fall short of human-level accuracy by over 30%. 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 stud-

ies, we present *eight key observations*, hoping that these insights will guide future improvements in VLMs for vision-centric tasks.

## 2 VLM2-Bench

As shown in Figure 2, **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

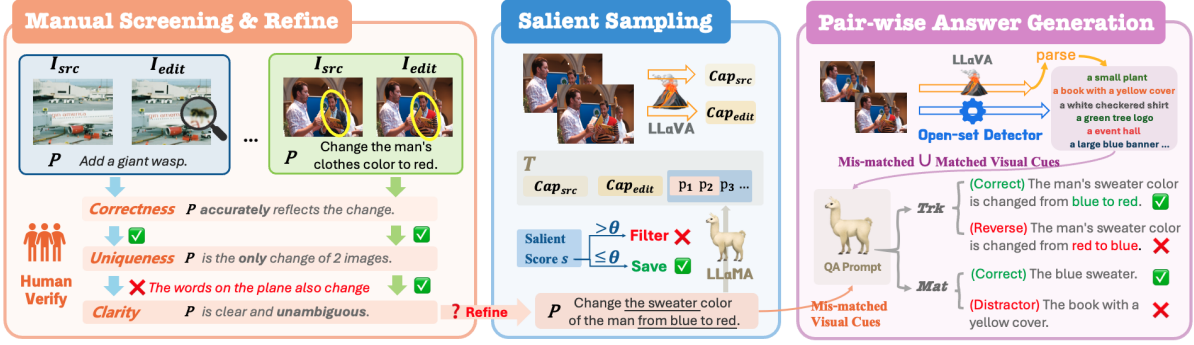


Figure 3: Construction of GC: (i) We start by manually verifying the edited image data based on three key criteria. (ii) A VLM is then prompted to generate captions for each image, followed by salient score-based filtering to retain the challenging cases. (iii) Finally, visual cues are extracted from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

only one of the two images and determine how it has changed. Rather than simply detecting a difference, 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 7 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:

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

where  $\mathcal{P} = \{p_1, p_2, \dots, p_{|\mathcal{P}|}\}$  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  $\theta$  (-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. Meanwhile, 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 visual consistency or change. Notably, we observe 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 actual capabilities. To mitigate this, we introduce consistency-pair validation, where for each statement (e.g., “X is Y”, with the answer being T), we generate a corresponding negation (e.g., “X is not Y”, with the answer being F). The model is only considered correct if it correctly answers both statements, ensuring consistency in its decision-making. (ii) **Counting (Cnt)** involves identifying the number of unique objects, requiring the model 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 instances of the same object, building on precise cue matching across multiple images.

**Data Collection.** We manually collect various categories of everyday objects (e.g., pets, cups) from multiple online resource<sup>2</sup>. 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  $\mathcal{I}_{-O_i}$ , 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)

PC aims to evaluate a model’s ability to link *person-centric 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.

<sup>2</sup><https://www.amazon.com/>, <https://lens.google/>, and <https://jellycat.com/>.

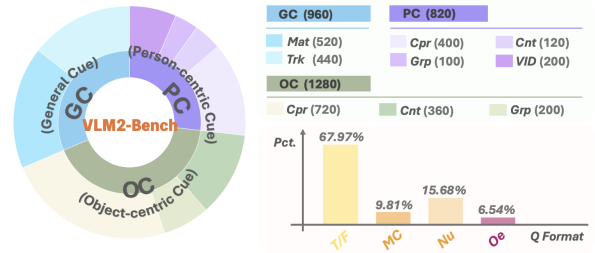


Figure 4: Statistical overview of **VLM2-Bench**. The pie chart shows the distribution of 9 subtasks across the 3 main categories of visual cues. The bar plot illustrates the percentage breakdown by question format.

**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  $I_{-i} \notin \mathcal{I}_{\mathcal{P}_i}$  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  $V_{-\mathcal{P}_i}$  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.

### 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



Baselines or Models	GC		OC			PC				Overall*	
	Mat	Trk	Cpr	Cnt	Grp	Cpr	Cnt	Grp	VID	Avg	$\Delta_{human}$
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	91.00	97.08	92.87	91.17	100.00	94.44	0.00
LLaVA-OneVision-7B	16.60	13.70	47.22	56.17	27.50	62.00	46.67	37.00	47.25	38.36	-56.08
LLaVA-Video-7B	18.53	12.79	54.72	62.47	28.50	62.00	66.91	25.00	59.00	41.37	-53.07
LongVA-7B	14.29	19.18	26.67	42.53	18.50	21.50	38.90	18.00	3.75	24.95	-69.49
mPLUG-Owl3-7B	17.37	18.26	49.17	62.97	31.00	63.50	58.86	26.00	13.50	40.89	-53.55
Qwen2-VL-7B	27.80	19.18	68.06	45.99	35.00	61.50	58.59	49.00	16.25	45.64	-48.80
Qwen2.5-VL-7B	35.91	43.38	71.39	41.72	47.50	80.00	57.98	69.00	46.50	55.86	-38.58
InternVL2.5-8B	21.24	26.03	53.33	55.23	46.50	51.50	60.00	52.00	5.25	45.73	-48.71
InternVL2.5-26B	30.50	30.59	43.33	51.48	52.50	59.50	59.70	61.00	21.75	48.58	-45.86
Gemini-2.0-flash	1.54	14.61	51.67	35.57	23.00	49.00	30.24	21.00	-	28.33	-66.11
Claude-3.7-sonnet	33.72	36.41	74.44	73.02	64.50	67.50	67.00	60.00	61.25	59.57	-34.87
GPT-4o-2024-08-06	37.45	39.27	74.17	80.62	57.50	50.00	90.50	47.00	66.75	59.56	-34.88
GPT-4o-2024-11-20	18.53	29.68	81.67	77.08	57.50	56.00	78.39	47.00	76.55	55.73	-38.71

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, multi-choice (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

#### 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:

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

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:

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

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  $\alpha$  as an error amplification factor. The final accuracy over  $n$  cases is:

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

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

**Open-ended** (*Video Identity Describing*): We use GPT-4o 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-4o’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 open-source VLMs that support multiple-image or video input: LLaVA-OneVision (Li et al.,

Model	Matching ( <i>Mat</i> )				Tracking ( <i>Trk</i> )			
	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-4o	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.

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 models GPT-4o (Hurst et al., 2024), Claude-3.7-sonnet, and Gemini-2.0-flash for comparison for comparison.

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

### 3.3 Results and Findings

**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 VLM<sup>2</sup>-Bench, covering a total of nine subtasks.

**Finding I: Simple tasks for humans pose significant challenges for VLMs.** We observe that humans achieve near-perfect accuracy across most tasks in our VLM<sup>2</sup>-Bench. In contrast, even state-of-the-art closed-source models perform significantly lower than humans. 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. 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 context before they can determine that the swapped cue has been transformed. This task requires a

Res.	GC		OC			PC		
	<i>Mat</i>	<i>Trk</i>	<i>Cpr</i>	<i>Cnt</i>	<i>Grp</i>	<i>Cpr</i>	<i>Cnt</i>	<i>Grp</i>
<b>Qwen2.5-VL-7B</b>								
Origin	35.91	43.38	71.39	41.72	47.50	80.00	57.98	69.00
↓ ×2	25.10	40.18	64.17	45.75	42.50	76.00	60.45	70.00
↓ ×4	19.69	33.33	52.78	42.25	33.00	64.50	57.15	61.00
↓ ×8	13.90	24.66	43.33	43.22	24.00	57.00	48.65	52.00
↓ ×16	9.27	18.72	34.17	38.86	22.50	45.50	47.01	41.00
<b>InternVL2.5-8B</b>								
Origin	21.24	26.03	53.33	55.23	47.50	51.50	60.00	52.00
↓ ×2	10.42	19.63	72.50	53.33	45.00	50.50	53.67	50.00
↓ ×4	11.97	20.09	69.72	52.77	47.00	51.00	54.99	49.00
↓ ×8	10.04	16.89	68.33	49.96	45.00	52.00	52.86	49.00
↓ ×16	3.47	14.16	61.39	43.30	47.50	49.00	51.06	50.00

Table 3: Models’ performance at the original resolution and with various compression levels. Results show a clear performance decline as image quality decreases, indicating that VLM<sup>2</sup>-Bench requires models to perceive and distinguish fine-grained visual details.

deeper understanding of how cues relate to each other across different instances. In contrast, *Trk* challenges models with mismatched cues due to add/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 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.

### 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.

### 3.4 Visual Bias Sanity Check

To assess whether models genuinely rely on fine-grained visual cues—rather than shortcut biases such as global layout or coarse semantics, we conduct a sanity check via image resolution ablation. Specifically, we evaluate two models (Qwen2.5-VL-7B and InternVL2.5-8B) under different levels of image compression, reducing the resolution by factors of 2, 4, 8, and 16.

As shown in Table 3, we observe a consistent performance drop as the image resolution decreases. This trend suggests that models do rely on detailed visual cues to perform well, rather than exploiting high-level layout or textual artifacts. These results highlight two key insights: (i) Our benchmark tasks indeed require models to perceive and distinguish fine-grained visual differences, rather than exploiting shallow biases. (ii) The performance sensitivity to visual degradation provides evidence that top-performing models are engaging in genuine visual understanding—which supports the benchmark’s role in probing visual linking ability under realistic, perception-driven settings.

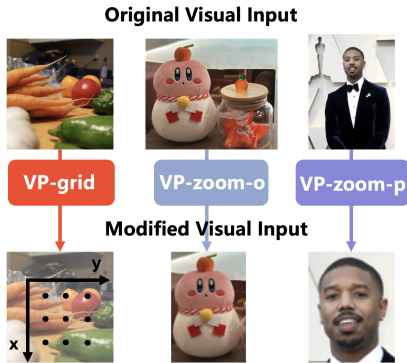
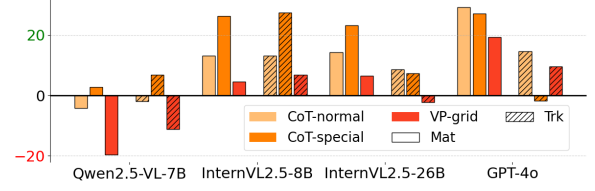


Figure 5: Visualization of three visual prompting (VP) approaches we adopted in Section 4. From left to right: the VPs used for GC, OC, and PC, respectively.

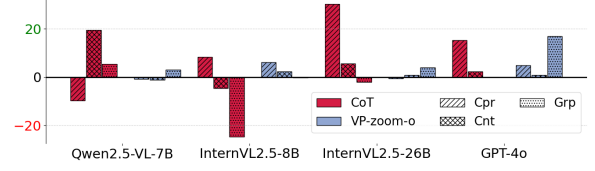
## 4 How Prompting Methods Affect VLMs

In this section<sup>3</sup>, 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-4o, and explore different approaches of CoT (Kojima et al., 2022; Wei et al., 2023) and visual prompting (VP) (Lei et al., 2024;

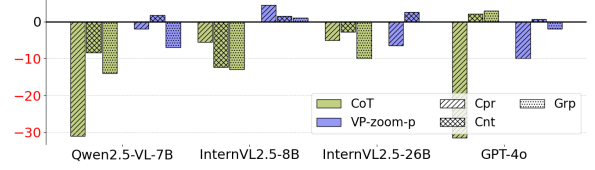
<sup>3</sup>Due to space limits, we reference most case studies, figures, and details in the Appendix within this section.



(a) Results of CoT-normal, CoT-special, and VP-grid on GC.



(b) Results of CoT and VP-zoom-o on OC.



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

Figure 6: Performance gains and losses (%) when applying different prompting methods on VLM2-Bench. (a) shows results on GC using CoT-normal, CoT-special, and VP-grid; (b) presents results on OC with CoT and VP-zoom-o; and (c) reports results on PC with CoT and VP-zoom-p. Detailed analyses are provided in Section 4.1, Section 4.2, and Section 4.3, respectively.

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 22) encourages the model to solve the task step by step, allowing it to reason through the problem. (ii) **CoT-special** (Table 23) guides the model to solve the task using a thought process closer to how humans typically approach it. (iii) **VP-grid** (Figure 12) is adapted from previous work (Lei et al., 2024) for our tasks, overlaying a dot matrix on the image as visual anchors to provide positional references and enhance the model’s performance.

**Finding IV: Reasoning in language aids models in logically linking visual cues.** From Figure 6a, 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 15, 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 depends on models’ ability to interpret both prompting cues and the visual content.** As shown in Figure 6a, VP-grid negatively impacts GC performance for QwenVL2.5, causing a significant drop compared to the vanilla approach. Figure 16 reveals that this decline stems from the model’s difficulty in interpreting the visual coordinates within the prompt, leading to misinterpretation of the cues and causing it to fail cases it originally answered correctly under the vanilla setting. However, as shown in Figure 17, GPT-4o successfully resolves a previously incorrect case by effectively leveraging the cues introduced through visual prompting while utilizing its strong visual perception abilities.

## 4.2 Probing for Object-centric Cue (OC)

**Methods.** (i) CoT (Table 22), and (ii) **VP-zoom-o** (Figure 13), which uses an open-set detector (Ren et al., 2024) to obtain bounding boxes. These boxes are then cropped to focus the model’s attention on object-centric cues. By eliminating irrelevant non-object cues and emphasizing the object-centric cues, this approach 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 6b, 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 18) 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 6b, we observe that for models with strong vision capabilities like GPT-4o, our VP-zoom-o method further enhances performance. For other models, this method at least en-

sures 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 22). (ii) **VP-zoom-p** (Figure 14) 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 person-centric cues, leading to a performance drop.** From Figure 6c, 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 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

**Advancements in vision-language models** have significantly broadened their capabilities (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). 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; He et al., 2025), 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 (Weerakoon et al., 2024; Yang et al., 2024; Jiang et al., 2024; Ye et al., 2025), 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 a critical role in guiding their future develop-

ment (Liang et al., 2024a; Yin et al., 2023; Chen et al., 2024a). These benchmarks typically focus on assessing the models’ fine-grained perception (Li et al., 2024a; Tong et al., 2024), reasoning abilities (Lu et al., 2022; Yu et al., 2023; Huang et al., 2025), commonsense knowledge (Yue et al., 2024; Wu et al., 2024), social intelligence (Li et al., 2025b), and robustness to input variations (Fan et al., 2025). In addition, evaluations targeting multi-image and 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; Li et al., 2025d), comparison (Wu et al., 2025; Jiao et al., 2024), and temporal reasoning (Liu et al., 2024b). However, existing benchmarks focus on evaluating VLMs’ ability to interpret visual cues based on their knowledge. In contrast, humans typically solve such tasks by explicitly matching visual 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-4o falls 34.70% 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.

## Limitations

VLM2-Bench focuses on evaluating visual cue linking but does not cover all possible scenarios. Additionally, while it provides valuable insights, its scale is limited, and model performance may not fully generalize to all real-world settings. Automated evaluation constraints limit the inclusion of open-ended questions in our benchmark, impacting the assessment of models’ vision-centric reasoning abilities. Expanding task diversity and refining evaluation methods (e.g., switching the one-shot evaluation scenario to multi-turn conversations (Li et al., 2025c)) remain important directions for future work. In future research, model self-play (Li et al., 2025a), self-correction (He et al., 2024), and synthetic data pretraining (Qin et al., 2025) may also be interesting to explore.

## 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*.
- Zhiyuan Fan, Yumeng Wang, Sandeep Polisetty, and Yi R. Fung. 2025. [Unveiling the lack of lvlm robustness to fundamental visual variations: Why and path forward](#). *Preprint*, arXiv:2504.16727.
- 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.
- Jiayi He, Hehai Lin, Qingyun Wang, Yi Fung, and Heng Ji. 2024. [Self-correction is more than refinement: A learning framework for visual and language reasoning tasks](#). *Preprint*, arXiv:2410.04055.
- Zhitao He, Sandeep Polisetty, Zhiyuan Fan, Yuchen Huang, Shujin Wu, and Yi R. Fung. 2025. [Mmboundary: Advancing mllm knowledge boundary awareness through reasoning step confidence calibration](#). *Preprint*, arXiv:2505.23224.
- Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. 2021. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*.
- Kung-Hsiang Huang, Hou Pong Chan, May Fung, Haoyi Qiu, Mingyang Zhou, Shafiq Joty, Shih-Fu Chang, and Heng Ji. 2025. [From pixels to insights: A survey on automatic chart understanding in the era of large foundation models](#). *IEEE Transactions on Knowledge and Data Engineering*, 37(5):2550–2568.
- 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 Wenhui 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 vision-language 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. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Cheng Li, May Fung, Qingyun Wang, Chi Han, Manling Li, Jindong Wang, and Heng Ji. 2025a. [Mentalarena: Self-play training of language models for diagnosis and treatment of mental health disorders](#). *Preprint*, arXiv:2410.06845.
- Hengzhi Li, Megan Tjandrasuwita, Yi R. Fung, Armando Solar-Lezama, and Paul Pu Liang. 2025b. [Mimeqa: Towards socially-intelligent nonverbal foundation models](#). *Preprint*, arXiv:2502.16671.
- Li Li, Peilin Cai, Ryan A Rossi, Franck Dernoncourt, Branislav Kveton, Junda Wu, Tong Yu, Linxin Song, Tiankai Yang, Yuehan Qin, et al. 2025c. A personalized conversational benchmark: Towards simulating personalized conversations. *arXiv preprint arXiv:2505.14106*.



- You Li, Heyu Huang, Chi Chen, Kaiyu Huang, Chao Huang, Zonghao Guo, Zhiyuan Liu, Jinan Xu, Yuhua Li, Ruixuan Li, et al. 2025d. Migician: Revealing the magic of free-form multi-image grounding in multimodal large language models. *arXiv preprint arXiv:2501.05767*.
- Paul Pu Liang, Akshay Goindani, Talha Chafekar, Leena Mathur, Haoifei Yu, Ruslan Salakhutdinov, and Louis-Philippe Morency. 2024a. Hemm: Holistic evaluation of multimodal foundation models. *arXiv preprint arXiv:2407.03418*.
- Paul Pu Liang, Amir Zadeh, and Louis-Philippe Morency. 2024b. Foundations & trends in multimodal machine learning: Principles, challenges, and open questions. *ACM Computing Surveys*, 56(10):1–42.
- 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 preprint 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*.
- Zeyu Qin, Qingxiu Dong, Xingxing Zhang, Li Dong, Xiaolong Huang, Ziyi Yang, Mahmoud Khademi, Dongdong Zhang, Hany Hassan Awadalla, Yi R. Fung, Weizhu Chen, Minhao Cheng, and Furu Wei. 2025. [Scaling laws of synthetic data for language models](#). *Preprint*, arXiv:2503.19551.
- 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 open-world models for diverse visual tasks](#). *Preprint*, arXiv:2401.14159.
- Jongyoon Song, Sangwon Yu, and Sungroh Yoon. 2024. Large language models are skeptics: False negative problem of input-conflicting hallucination. *arXiv preprint arXiv:2406.13929*.
- Qwen Team. 2025. [Qwen2.5-vl](#).
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. Fever: a large-scale dataset for fact extraction and verification. *arXiv preprint arXiv:1803.05355*.
- Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. 2024. Eyes wide shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9568–9578.
- Anne Treisman and Garry A. Gelade. 1980. [A feature-integration theory of attention](#). *Cognitive Psychology*, 12:97–136.
- Fei Wang, Xingyu Fu, James Y Huang, Zekun Li, Qin Liu, Xiaogeng Liu, Mingyu Derek Ma, Nan Xu, Wenxuan Zhou, Kai Zhang, et al. 2024a. Muirbench: A comprehensive benchmark for robust multi-image understanding. *arXiv preprint arXiv:2406.09411*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin 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 perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- Yi Wang, Yinan He, Yizhuo Li, Kunchang Li, Jiashuo Yu, Xin Ma, Xinhao Li, Guo Chen, Xinyuan Chen, Yaohui Wang, Conghui He, Ping Luo, Ziwei Liu, Yali Wang, Limin Wang, and Yu Qiao. 2024c. [Internvid: A large-scale video-text dataset for multimodal understanding and generation](#). *Preprint*, arXiv:2307.06942.
- Kasun Weerakoon, Mohamed Elnoor, Gershom Seneviratne, Vignesh Rajagopal, Senthil Hariharan Arul, Jing Liang, Mohamed Khalid M Jaffar, and Dinesh Manocha. 2024. [Behav: Behavioral rule guided autonomy using vlms for robot navigation in outdoor scenes](#). *Preprint*, arXiv:2409.16484.
- Cong Wei, Zheyang Xiong, Weiming Ren, Xinrun Du, Ge Zhang, and Wenhui Chen. 2024. Omniedit: Building image editing generalist models through specialist supervision. *arXiv preprint arXiv:2411.07199*.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. [Chain-of-thought prompting elicits reasoning in large language models](#). *Preprint*, arXiv:2201.11903.
- Haoning Wu, Hanwei Zhu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Chunyi Li, Annan Wang, Wenxiu Sun, Qiong Yan, et al. 2025. Towards open-ended visual quality comparison. In *European Conference on Computer Vision*, pages 360–377. Springer.
- Jialian Wu, Jianfeng Wang, Zhengyuan Yang, Zhe Gan, Zicheng Liu, Junsong Yuan, and Lijuan Wang. 2022. Grit: A generative region-to-text transformer for object understanding. *arXiv preprint arXiv:2212.00280*.
- 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 Linguistics: EMNLP 2024*, pages 7715–7731, Miami, Florida, USA. Association for Computational Linguistics.
- Lingfeng Yang, Yueze Wang, Xiang Li, Xinlong Wang, and Jian Yang. 2023. [Fine-grained visual prompting](#). *Preprint*, arXiv:2306.04356.
- Zhutian Yang, Caelan Garrett, Dieter Fox, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. 2024. [Guiding long-horizon task and motion planning with vision language models](#). *Preprint*, arXiv:2410.02193.
- Dongyu Yao and Boheng Li. 2023. Dual-level interaction for domain adaptive semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, pages 4527–4536.
- Dongyu Yao, Jianshu Zhang, Ian G. Harris, and Marcel Carlsson. 2024. [Fuzzllm: A novel and universal fuzzing framework for proactively discovering jailbreak vulnerabilities in large language models](#). In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4485–4489.
- 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](#). *Preprint*, arXiv:2408.04840.
- Jiayi 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. *arXiv preprint arXiv:2410.02736*.
- Mang Ye, Xuankun Rong, Wenke Huang, Bo Du, Nenghai Yu, and Dacheng Tao. 2025. A survey of safety on large vision-language models: Attacks, defenses and evaluations. *arXiv preprint arXiv:2502.14881*.
- 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 arXiv:2306.13549*.
- 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 arXiv:2308.02490*.
- Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yuet-ing 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 9127–9134.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.
- Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkan Yang, Yuanhan Zhang, Ziyue Wang, Hao-ran 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 multi-image understanding in vision and language models: Perception, knowledge, reasoning, and multi-hop reasoning. *arXiv preprint arXiv:2406.12742*.

## A Appendix Outline

In the appendix, we provide:

- **Appendix B** provides details on the licensing terms and usage rights for our benchmark.
- **Appendix C** presents the statistical analysis of the VLM2-Bench.
- **Appendix D** details on how we obtain the chance-level and human-level baselines.
- **Appendix E** elaborates more details on the construction of the VLM2-Bench.

Category	T/F	MC	Nu	Oe	Total
GC	960	–	–	–	960
OC	720	200	360	–	1,280
PC	400	100	120	200	820
Total	2,080	300	480	200	3,060

Table 4: Overview of query distribution across the three categories of VLM2-Bench. T/F = True/False, MC = multiple-choice, Nu = numerical, Oe = open-ended.

- **Appendix F** provides a deeper dive into the various prompting techniques we use.
- **Appendix G** a detailed breakdown and analysis of failure and success examples regarding different prompting methods.

## B Licencing and Intended Use

Our VLM2-Bench is available under the CC-BY 4.0 license for academic use with proper attribution. The images, videos, and annotations in this benchmark are intended solely for research purposes. These data were sourced from publicly available online platforms, and while efforts were made to use them responsibly, explicit permissions may not have been obtained for all content. Users are responsible for ensuring that their use of the data complies with applicable intellectual property laws and ethical guidelines. We encourage users to verify the sources and ensure compliance with any terms of service or licensing agreements.

## C VLM2-Bench Statistics

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

### C.1 Overall Composition

Table 4 provides a detailed summary of the total query counts across different categories and subtasks in our benchmark. The dataset is structured into three primary categories: General Cue (GC),

Object-centric Cue (OC), and Person-centric Cue (PC), comprising a total of 3,060 visual-text query pairs.

The General Cue (GC) category consists of 960 queries, which include 260 Matching (Mat) true/false pairs, resulting in 520 queries, and 220 Tracking (Trk) true/false pairs, leading to 440 queries.

The Object-centric Cue (OC) category contains 1,280 queries, covering three subtasks: Comparison (Cpr) with 360 true/false pairs (720 queries), Counting (Cnt) with 360 numerical queries, and Grouping (Grp) with 200 multiple-choice questions.

Lastly, the Person-centric Cue (PC) category includes 820 queries, comprising 200 Comparison (Cpr) true/false pairs (400 queries), 120 Counting (Cnt) numerical queries, 100 Grouping (Grp) multiple-choice questions, and 200 Free-form (VID) open-ended queries.

Overall, these components collectively sum up to 3,060 visual-text query pairs, offering a comprehensive benchmark for evaluating vision-language models across various types of contextual cues.

### C.2 Details per Subtask and Question Type

#### General Cue (GC).

**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.

**Tracking (Trk).** We design 220 T/F pairs that test an understanding of object or entity continuity across frames. For example, a question might ask whether the same object reappears in subsequent frames. Each T/F pair similarly results in two queries, totaling 440.

**Object-centric Cue (OC).** All the visual query cases are built upon the 360 image sequences we construct. Details about image sequences can be found 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 (MC) questions that ask about grouping objects according to certain criteria (e.g., AAB, ABC, AAAB, AABC, ABCD). Each question presents multiple group-configuration options plus a “None” option, which can serve as either a correct or distractor choice. For image sequences of length 4, the options include various plausible groupings (two-of-a-kind, three-of-a-kind, etc.) along with at least one 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 numerical 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.

**Open-ended (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 evaluation of vision-language models. Figure 4 in the main paper illustrates the distribution of these subtasks and their question-format breakdown. We believe that the richness and diversity of VLM2-Bench make it a robust platform for advancing multimodal research.

## D Baselines

### D.1 Chance-level

In this part, we explain the calculation of chance-level accuracy for all subtasks in Table 1.

**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).

A question pair example is shown in Table 5.

#### 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: **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: **F**

Table 5: Example of True-False paired questions in GC-Mat, with a positive and negative question.

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.
- **Negative Question:** A negated version of the positive question, often by inserting "not" after the verb.

An example is shown in Table 6.

**Positive Question:**

"Given the images, the claim 'The pets in these images *are* the same pet.' is right."

GT Answer: **T**

**Negative Question:**

"Given the images, the claim 'The pets in these images *are not* the same pet.' is right."

GT Answer: **F**

Table 6: Example of True-False paired questions in OC-Cpr, with a positive and negative question.

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  $\{1, 2, \dots, L\}$ , 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:

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

where

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

and the weight is defined as

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

with  $L_{\max} = 4$  being the maximum sequence length in our dataset.

**OC-Cnt Task:** The OC-Cnt task exhibits the following distribution:

- Length 2: 80 sequences (22.2%)
- Length 3: 120 sequences (33.3%)
- Length 4: 160 sequences (44.4%)

Thus, the overall chance level accuracy is obtained as the weighted average:  $Acc_{OC-Cnt} = \frac{80 E(2) + 120 E(3) + 160 E(4)}{360} \approx 34.88\%$ .

**PC-Cnt Task:** For the PC-Cnt task, the sequence distribution is:

- Length 2: 30 sequences (25.0%)
- Length 3: 25 sequences (20.8%)
- Length 4: 65 sequences (54.2%)

Accordingly, the overall chance level accuracy is given by:  $Acc_{PC-Cnt} = \frac{30 E(2) + 25 E(3) + 65 E(4)}{120} \approx 34.87\%$ .

## D.2 Human-level

To facilitate human participants in providing responses to our questions, we integrated all model-prompted questions and answer choices into a graphical user interface (GUI), as illustrated in Figure 7. 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.

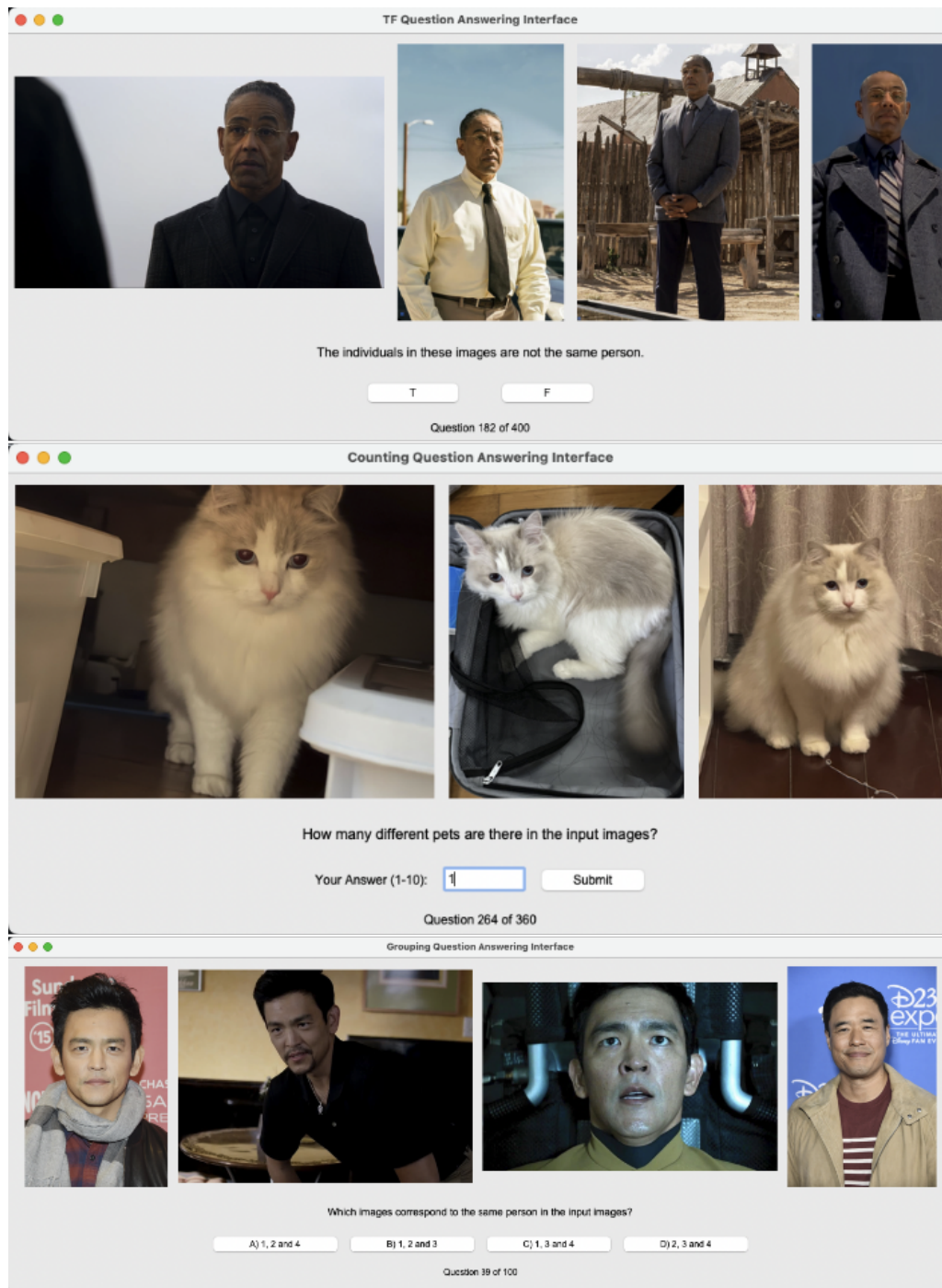


Figure 7: The GUI used for human-level testing.



## E More details on Benchmark Construction

### E.1 GC (General Cue)

**Manual Screening and Refine.** Figure 8 demonstrates the Graphic User Interface (GUI) we build for manually screening image editing data.

**Salient Sampling.** The pseudocode in Figure 9 and Table 7 displays the calculation process for the salient sampling score mentioned in Section 2.1.

**Prompts for Pair-wise Answer Generation.** Table 8 and 9 provides the complete prompts used to generate pair-wise answers for our evaluation tasks. The prompts were designed to instruct the language model to produce two distinct answers—a positive (T) answer and a negative (F) answer—for each task. The dual-answer format is intended to capture both the expected response and its direct opposite, thereby offering a more balanced insight into the model’s understanding.

### 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 10. In total, we manually collected 320 images for objects.

**Main Meta-Object Selection.** We predefine 8 types of common objects, with each type containing 5 meta-objects to ensure a class-balanced sampling and avoid long-tail distribution (Yao and Li, 2023). For each meta-object, we collect four images that represent the same object from different angles and scene conditions.

**Distractor Meta-Object Selection.** To build meaningful object image sequences, we introduce visually distractive elements for each main meta-object, referred to as “distractor meta-objects”. Specifically, for each main meta-object, we collect four additional images that belong to different but visually similar meta-objects within the same object category. These images are selected following predefined visual cue confusion principles, ensuring that they provide meaningful challenges for vision language models. We ensure that each distractor image belongs to a different distractor meta-object, fundamentally guaranteeing that the count of different meta-objects in the final constructed sequence strictly follows our design. The principle of selecting distractor meta-objects is illustrated in the outer ring of Figure 10.

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

- **Plush Objects:** Images of plush toys are entirely sourced from the [Jellycat website](#) and its review sections, where diverse user-uploaded images provide a wide variety of object angles and scenes.
- **Pet Objects:** For the pet category of meta-objects, 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 10.

**Images Sequence Construction.** The construction of image sequences in OC (a total of 360 sequences) follows the structure in Table 10. More specific details are listed below:

#### Two-Image Sequences (`image_seq_len = 2`)

1. **Main Meta-Object Only (AA):** Two images are randomly sampled from the same main meta-object. 40 sequences are constructed (one for each main meta-object).
2. **Main Meta-Object + Distractor Meta-Object (AB):** One image is randomly selected from the main meta-object, and one from the corresponding distractor meta-object. 40 sequences are constructed.

#### Three-Image Sequences (`image_seq_len = 3`)

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

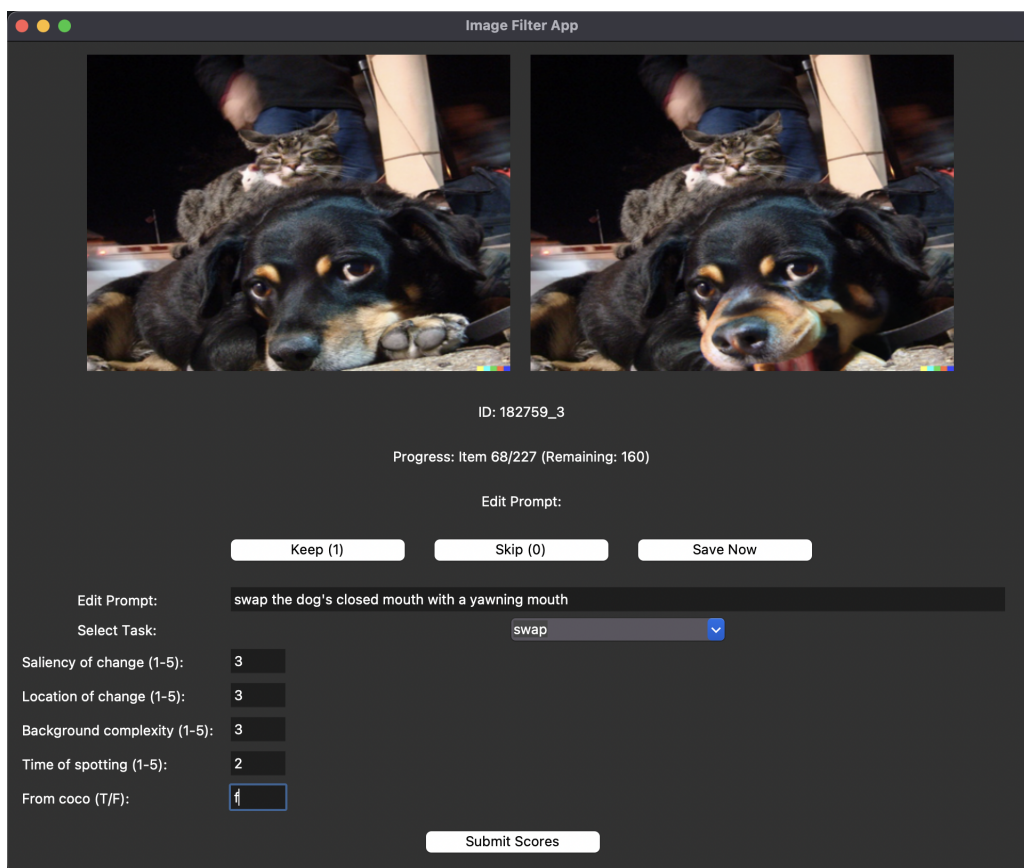


Figure 8: The GUI used for manually screening image editing data and refining edited prompts in General Cue (GC).

Supposed you are looking at two images:

Image 1: **<Cap\_src>**

Image 2: **<Cap\_edit>**

From Image 1 to Image 2, the change can be summarized as: **<P>**

Table 7: Template for salient-score calculation, which contain three placeholders for each sample.

### #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 mis-matched cue types, here shows the question for 'add' from the "Add/Remove" category). 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

*Editing Information:*

<Edit Prompt>

*Description:*

<Description>

Table 8: Prompt for generating paired answers in the Matching (Mat) subtask of General Cue (GC).

**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**

*Editing Information:*

<Edit Prompt>

*Description:*

<Description>

Table 9: Prompt for generating paired answers in the Tracking (Trk) subtask of General Cue (GC).



---

**Algorithm 1** Salient Score Computation

---

```
1 # cap_src: caption for the source image
2 # cap_edit: caption for the edited image
3 # T: template for constructing a paragraph
4 # P: editing prompt
5 input_text = concat(cap_src, cap_edit, T)
6 in_tokens = tokenizer.encode(input_text)
7 out_tokens = tokenizer.encode(P)
8 log_sum = 0
9 tokens = in_tokens
10
11 # Model Forward Pass
12 for i in range(1, len(out_tokens)):
13     outputs = model(tokens)
14     logits = outputs.logits
15
16     # Extract log probability of next token
17     probs = log_softmax(logits[0], -1, :])
18     prob = probs[out_tokens[i]]
19     log_sum += prob
20
21     # Update Input Sequence
22     tokens = concat(tokens, out_tokens[i])
23
24 # Normalize the total log probability as the salient_score
25 salient_score = log_sum / len(out_tokens)
26
27 # Return: salient_score
```

---

Figure 9: Pseudocode for salient score computation in the phrase of Salient Sampling in the construction of GC.

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.
3. **Main Meta-Object + Distractor Meta-Objects (AABC)**: Two images are selected from the main meta-object, while two are selected from different distractor meta-objects. 40 sequences are constructed.
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 sequences are constructed.

**Question Templates.** Table 11, 12 and 13 list detailed standard question templates (with format instructions) for the Object-centric Cue task, including 3 subtasks: Comparison (cpr), Counting (Cnt), and Grouping (Grp).

#### E.3 PC (Person-centric Cue)

**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



Figure 10: 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.	T	2	-
2	AB	1 image $I_i$ from $\mathcal{I}_{O_i}$ and 1 image $I_{-i}$ from distractor set $\mathcal{I}_{-O_i}$ , 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.	T	3	-
3	AAB	2 images from the same object $O_i$ , randomly sampled as $\mathcal{I}_{O_i} = \{I_i, I_j\}$ and 1 $I_{-i}$ from distractor set $\mathcal{I}_{-O_i}$ , 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 $\{I_{-i}, I_{-j}\}$ from distractor set $\mathcal{I}_{-O_i}$ , 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.	T	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 $I_{-i}$ from distractor set $\mathcal{I}_{-O_i}$ , 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 $\{I_{-i}, I_{-j}\}$ from distractor set $\mathcal{I}_{-O_i}$ , 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 $\{I_{-i}, I_{-j}, I_{-k}\}$ from distractor set $\mathcal{I}_{-O_i}$ , randomly shuffled.	F	3	$[\ ]$

Table 10: Summary of multi-images sequence construction for Object-centric Cue (OC) tasks.

#### OC-Cpr Positive Question:

*Judge the following statement based on the images: ‘The {obj}s in these images are the same {obj}.’ Provide only one correct answer: ‘T’ (True) or ‘F’ (False). Respond with either ‘T’ or ‘F’.*

GT Answer: **T**

#### OC-Cpr Negative Question:

*Judge the following statement based on the images: ‘The {obj}s in these images are **not** the same {obj}.’ Provide only one correct answer: ‘T’ (True) or ‘F’ (False). Respond with either ‘T’ or ‘F’.*

GT Answer: **F**

#### 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: **3** (Example Answer)

Table 12: The question template used for the counting (Cnt) subtask of Object-centric Cue (OC).

Table 11: Question templates used for consistency-pair evaluation in the Comparison (Cpr) subtask of Object-centric Cue (OC).

**OC-Grp Question:**

*Answer the following question based on this rule: You only need to provide \*ONE\* correct answer, selecting from the options listed below. For example, if you think the correct answer is 'B) 1 and 2', your response should 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 13: The question template used for the grouping (Grp) subtask of Object-centric Cue (OC).

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 introduce challenging distractors in our sequences, we compute the CLIP embedding for every image and store these embeddings in a reference base. When a distractor image is needed, we perform an image-to-image similarity search within this base to identify the most visually similar image that originates from a different meta-human. This fine-grained matching ensures that the distractor image closely resembles the main meta-human’s image, leading to more challenging image sequences.

**Discussion on Why Objects Require Dedicated Distractors, While Humans Do Not.** In object-centric tasks, objects are categorized into eight distinct types, with substantial differences among different types (e.g. pets and bags). Therefore, each main meta-object requires dedicated distractors from the same object type to ensure meaningful comparisons. In contrast, humans belong to a single category, meaning that any meta-human can serve as a distractor for another. Given that we compute CLIP embeddings to select visually similar distractors, the constructed image sequences already present a significant challenge without the need for type-specific distractors. We also ensure diversity by selecting five main meta-humans for each race-gender pair, providing a sufficiently large

pool from which to choose suitable distractors. Corresponding to our hypothesis, in the final curated sequences, most distractor meta-humans chosen were of the same race or gender as the main meta-human. Additionally, as shown in Table 1, these curated image sequences along with our designed questions effectively challenge tested models, revealing their limited performances in visually linking matching cues on person-centric data.

**Images Sequence Construction.** The construction of image sequences in PC (a total of 260 sequences) follows the structure in Table 14. More specific details are listed below:

#### **Two-Image Sequences (image\_seq\_len = 2)**

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

#### **Three-Image Sequences (image\_seq\_len = 3)**

1. **Main Meta-Human Only (PPP):** Three images are randomly sampled from the same main meta-human. 20 sequences are constructed.
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.
3. **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. 30 sequences are constructed.
2. **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. 20 sequences are constructed.

3. **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. 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. 30 sequences are constructed.

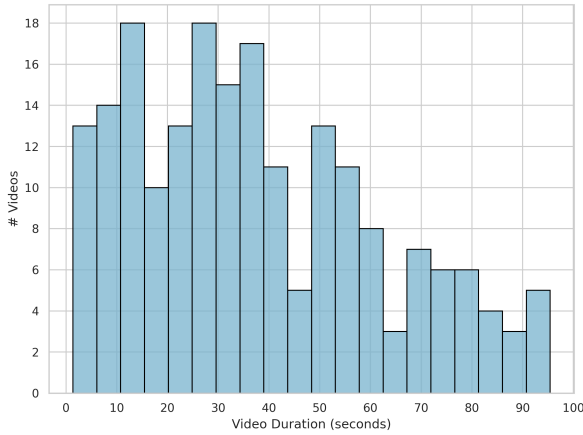


Figure 11: Distribution of video duration in the subtask of Video Identity Description (VID) in PC.

**Video Construction.** The video data for this benchmark is manually collected from Shutterstock<sup>4</sup>. 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 **0-100s** time range. Figure 11 displays the sketch of concatenated video length distribution. The final compositions followed two formats:

**$P \rightarrow P$  format** : A direct concatenation of two distinct clips (same length for each clip).

**$P \rightarrow \neg P \rightarrow P$  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).

Regardless of the different default sampling methods for our baseline models in Table 15, both  $P \rightarrow \neg P$  and  $P \rightarrow \neg P \rightarrow P$  formats ensure that every video clip has frames included while sampling:

- **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.

Model Name	Uni	FPS
LLaVA-OneVision-7B	✓	✗
LLaVA-Video-7B	✓	✗
LongVA-7B	✓	✗
mPLUG-Owl3-7B	✓	✗
Qwen2-VL-7B	✗	✓
Qwen2.5-VL-7B	✗	✓
InternVL2.5-8B	✓	✗
InternVL2.5-26B	✓	✗
GPT-4o	✓	✗
Claude-3.7-sonnet	✓	✗

Table 15: Comparison of different video sampling methods of VLMs, including Uniform Sampling (Uni) and FPS Sampling (1fps).

Thus, by maintaining the integrity of each clip’s temporal structure, both  $P \rightarrow \neg P$  and  $P \rightarrow \neg P \rightarrow P$  formats effectively ensure that every clip contributes frames to the final sampled frame input for all models.

**Question Templates.** Table 16, Table 17, Table 18, and Table 19 present the detailed standard question templates for the Person-centric Cue task, covering PC-Cpr, PC-Cnt, PC-Grp, and PC-VID.

<sup>4</sup><https://www.shutterstock.com>

Num	Src	Process of Image Sequences Construction	Cpr	cnt	Grp
2	<b>PP</b>	2 images from the same person $P_i$ , randomly sampled as $\mathcal{I}_{P_i} = \{I_i, I_j\}$ , and shuffled.	T	2	-
2	<b>PQ</b>	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	<b>PPP</b>	3 images from the same person $P_i$ , randomly sampled as $\mathcal{I}_{P_i} = \{I_i, I_j, I_k\}$ , and shuffled.	T	3	-
3	<b>PPQ</b>	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	<b>PQR</b>	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	$[\ ]$
4	<b>PPPP</b>	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.	T	4	-
4	<b>PPPQ</b>	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	<b>PQQR</b>	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	<b>PQRV</b>	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 14: Summary of multi-images sequence construction for Person-centric Cue (PC) tasks.

**PC-Cpr Positive Question:**

*Judge the following statement based on the images: 'The individuals in these images are the same person.' Provide only one correct answer: 'T' (True) or 'F' (False). Respond with either 'T' or 'F'.*

GT Answer: **T**

**PC-Cpr Negative Question:**

*Judge the following statement based on the images: 'The individuals in these images are **not** the same person.' Provide only one correct answer: 'T' (True) or 'F' (False). Respond with either 'T' or 'F'.*

GT Answer: **F**

Table 16: Question templates used for consistency-pair evaluation in the Comparison (Cpr) subtask of Person-centric Cue (PC).

**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 17: The question template used for the counting (Cnt) subtask of Person-centric Cue (PC).

## 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 an evaluator model (Yao et al., 2024) (here GPT-4o) using the scoring prompts detailed in Tables 20 and 21. Specifically, for videos following a  $\mathcal{P} \rightarrow \neg\mathcal{P}$  sequence, GPT-4o 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  $\mathcal{P} \rightarrow \neg\mathcal{P} \rightarrow \mathcal{P}$  (PQP) pattern, the evaluation is more nuanced. The evaluator model (GPT-4o) 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 22. 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 23 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 12 displays a complete version of Visual Prompting with Grid assistance (VP-grid). Here we follow (Lei et al., 2024) to print a set of dot matrix onto the input image, accompanied by the image order dimension concatenated with Cartesian coordinates as (*image order index*, *column index*), *row index*). In the detailed textual prompt design, we also integrated references and explanations for the grids, allowing VLMs to

**PC-Grp Question:**

*Answer the following question according to this rule: You only need to provide \*ONE\* correct answer, selecting from the options listed below. 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 18: The question template used for the grouping (Grp) subtask of Person-centric Cue (PC).

**PC-VID Question:**

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

Table 19: The question template used for the Video Identity Description (VID) subtask of Person-centric Cue (PC).

leverage this visual assistance as spatial and visual matching references.

**VP-zoom-o for OC.** In Figure 13, 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 process is similar to that of OC (Figure 14). We use a face detection model (Geitgey, 2016) to “zoom in” on the individual’s face and occlude other irrelevant information.

## G Case Study

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 InternVL2.5-8B’s performance by over 25% than the standard query in both Matching and Tracking tasks for General Cue. While for the traditional CoT-normal prompting technique, this boost is only 13%. The CoT-special prompt (Table 23) directs the model through four explicit steps: understanding the question, perceiving (listing elements), connecting (comparing and reasoning), and concluding. This structured approach mirrors the human process of visual matching and is effective even for

a rather smaller model like InternVL2.5-8B, which might otherwise struggle with the ambiguity of a complex generic step-by-step instruction (which we will discuss later in the next Subsection G.2).

For example, in the provided InternVL2.5-8B response Figure 15, 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 only ensures that all pertinent visual cues are captured and processed but also reduces errors by structuring the logical flow of reasoning. The CoT-special prompt’s explicit instructions help InternVL2.5-8B align visual information with textual descriptions more effectively, thus enhancing overall performance. Compared to the less specific CoT-normal prompt—which may leave the model with gaps in reasoning—the CoT-special prompt provides clear, task-specific guidance that is essential for complex visual reasoning tasks, as evidenced by the substantial performance improvement.



**#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 20: 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 21: 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 22: CoT prompt for GC (here we denote as CoT-normal to distinguish it from the CoT-special in Table 23 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)

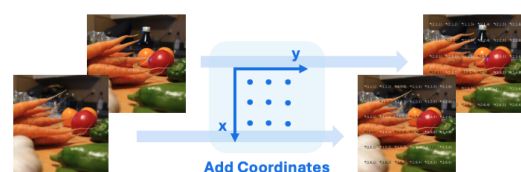
- List every details in each image respectively.
- Note positions and attributes of elements.

#### 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 23: 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 12: An illustration of how VP-grid works for GC.

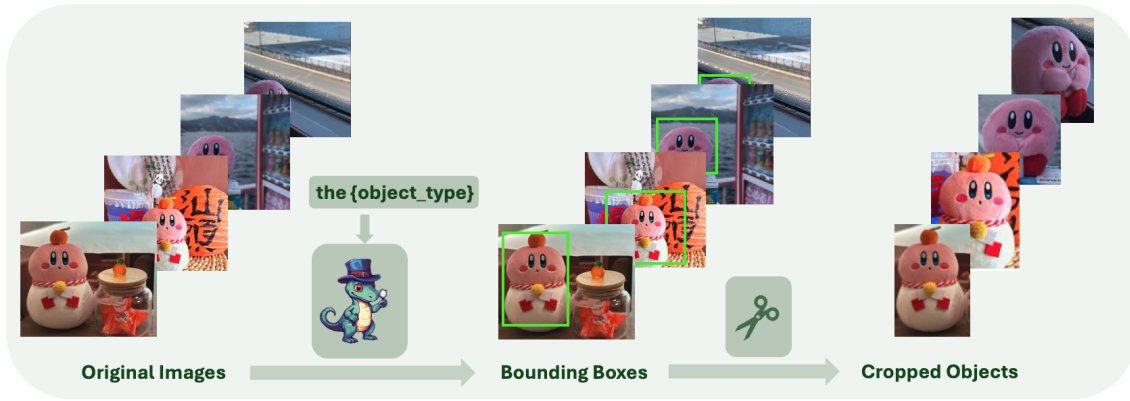


Figure 13: An illustration of how VP-zoom-o works for OC.

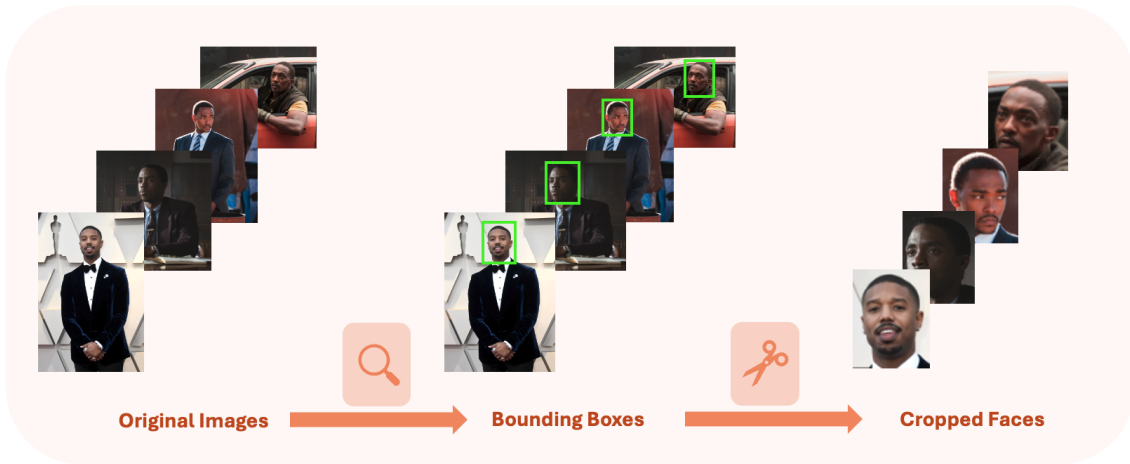


Figure 14: An illustration of how VP-zoom-p works for PC.

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

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.

**An example failure case in VP-grid.** However, this approach does not yield consistent improvements across all models. For instance, the Qwen2.5-VL-7B model demonstrates a significant performance drop—nearly 20%—when using VP-grid. An example failure case is in Figure 16. Our analysis reveals that although the model cor-

rectly identifies visual elements (e.g., a pedestrian with a high-visibility vest at coordinates  $(2, 5, 3)$ ), it fails to properly interpret the image sequence. Specifically, the model incorrectly associates the coordinates  $(2, 5, 3)$  with the first image, rather than the second, despite the explicit definition provided in the textual prompt. This misinterpretation leads to erroneous linking of visual matching cues and subsequent faulty reasoning. We suspect that the underlying issue is the limited semantic comprehension capability of the relatively smaller 7B model, which struggles with complex, predefined spatial instructions and visual assistance.

**An example of success case in VP-grid.** In contrast to models that often misinterpret or neglect spatial cues provided by VP-grid—leading to errors such as mismatching image indices—GPT-4o successfully leverages these visual prompts to achieve correct visual-textual alignment. In the example at Figure 17, the model identifies the cat’s nose at coordinates  $(1, 2, 4)$  in the first image and at  $(2, 2, 4)$



Image1



Image2

**Question:** Is the answer 'The vase on top of the table was removed.' correct for the given question: 'What key visual change can be observed from the first image to the second image?'?

**GT answer:** T **Vanilla Response:** F

**CoT-Special Response:** Step 1. Understand the Question ...

Step 2. Perceive (List Elements)

Image 1: Vase with flowers on the table; Chandelier above...

Image 2: No vase on the table; Chandelier above...

Step 3. Connect (Compare & Reason)

Both images show a {chandelier, kitchen area, bowl of fruit, window...}

The key visual change is the absence of the vase on the table in Image 2.

Step 4. Conclude (Answer the Question)

The statement "The vase on top of the table was removed" is correct for the key visual change observed between the two images.

MY\_ANSWER: T

Figure 15: Case study illustrating how CoT-special improves performance of the subtask of Tracking (Trk) in General Cue (GC). The model, InternVL2.5-8B, demonstrates a step-by-step reasoning process: In Step 2, it identifies key 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 compares the images, recognizing that while many elements remain unchanged (e.g., chandelier, kitchen area, fruit bowl, window), the vase's removal is the primary difference. In Step 4, the model concludes that the statement "The vase on top of the table was removed" accurately reflects the visual change, leading to the correct answer.



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-4o 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-4o with VP-grid can be attributed to its larger model capacity. Although the detailed architecture of GPT-4o 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-4o 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

guidance improves performance rather than causing confusion.

The success of GPT-4o 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-4o 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%.

### G.3 Case for CoT prompting in 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 6b, 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 18 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 CoT 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

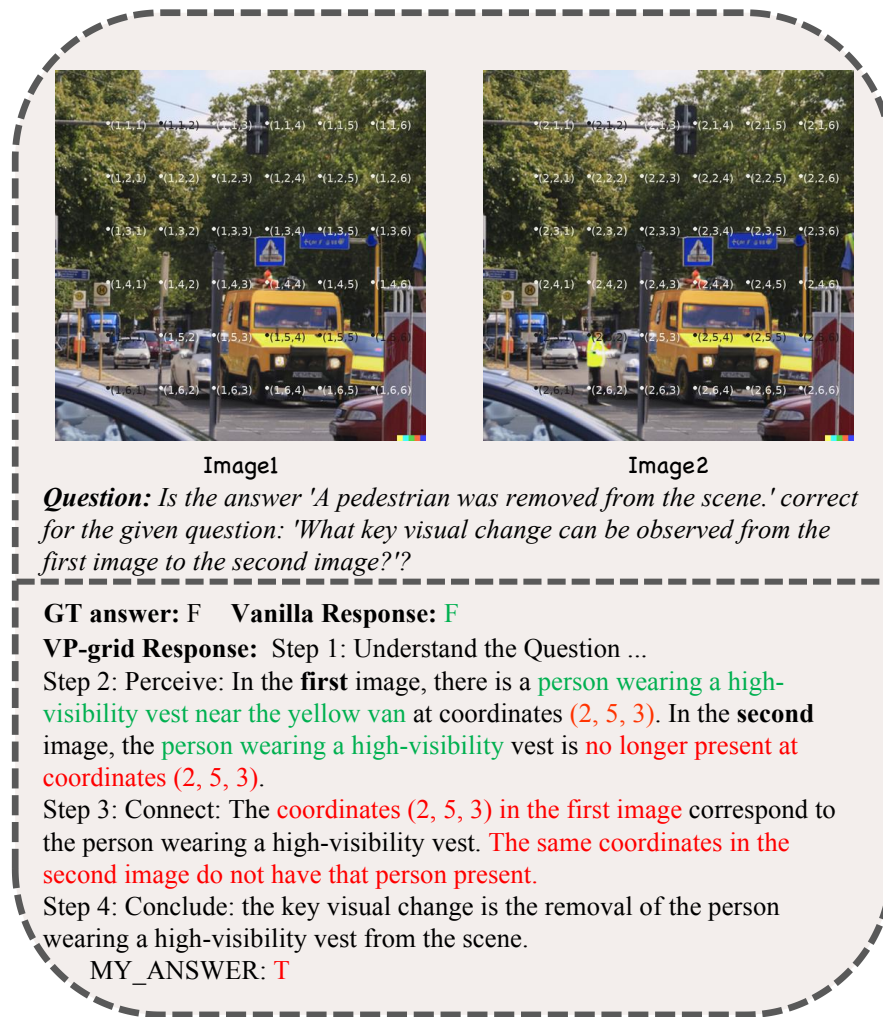


Figure 16: Case study illustrating why VP-grid leads to performance degradation in Qwen2.5-VL-7B. The model correctly identifies visual elements (e.g., a pedestrian with a high-visibility vest at coordinates (2, 5, 3)) but fails to interpret the image sequence correctly. It mistakenly associates the coordinates with the first image instead of the second, despite the explicit definition in the textual prompt. This misinterpretation results in incorrect visual cue linking and faulty reasoning, highlighting the model’s difficulty in handling structured spatial instructions and visual prompts.

- 3, where the dinosaur holds a keyboard rather than a skateboard.
2. **Variability in Descriptive Language:** The open-ended language generated by the CoT approach can lead to inconsistent descriptions. In this case, the model generalized the visual cue of a "dinosaur design" without capturing the specific attribute (i.e., the object the dinosaur is holding), which is crucial for correct grouping.

Thus, the lack of structured guidance in the CoT prompt leads to the dropping or misinterpretation of critical cues, resulting in incorrect grouping decisions for multi-image sequences in OC tasks. This analysis underscores the importance of more de-

tailed structured intermediate reasoning strategies, such as those provided by a tailored CoT-special prompt, to ensure that all relevant visual details are captured and compared accurately.



**Question:** Is the answer 'The cat's pink nose was turned black' correct for the given question: 'What key visual change can be observed from the first image to the second image?'?

**GT answer:** T **Vanilla Response:** F

**VL-grid Response:** (Understanding question context...) ... focusing on the cat's nose color change.

(Perception step...) Identify cat's nose location:

Image 1: Cat's white/pink nose at coordinates (1, 2, 4)

Image 2: Cat's black nose at coordinates (2, 2, 4)

(Connecting step...)

Compare cat's nose color at the same coordinates in both images:

**Image 1:** Nose appears lighter (likely pink).

**Image 2:** Nose appears darker (black).

(Concluding step...) The described change: "The cat's pink nose was turned black," is indeed the key visual change.

MY\_ANSWER: T

Figure 17: Case study demonstrating why VP-grid leads to performance improvement for GPT-4o. Unlike models that often misinterpret or overlook spatial cues, GPT-4o effectively uses VP-grid to align visual and textual information. In the example shown in Figure 17, the model correctly identifies the cat's nose at coordinates (1, 2, 4) in the first image and (2, 2, 4) in the second, accurately capturing the visual change in the attribute (from a lighter pink to a darker black). This success highlights GPT-4o's ability to handle structured spatial prompts and improve performance through visual prompting.




Image1




Image2




Image3

**Question:** Which images show the same bag in the input images?  
**Choices:** A) 1 and 3; B) 1 and 2; C) 2 and 3; D) None

**GT answer:** D) None      **Vanilla Response:** D) None

**CoT Response:**  
 The first image shows a backpack with a green dinosaur design worn by a child standing in a classroom. The second image shows a close-up of a backpack with a similar green dinosaur design, including sunglasses, but without the context of a person wearing it. The third image shows another backpack with a design featuring a dinosaur with sunglasses. The second and third images both have dinosaurs wearing sunglasses, which suggests they are the same bag.

**MY\_ANSWER:** C) 2 and 3

Figure 18: Case study illustrating why CoT leads to performance degradation. In the example shown in Figure 18, InternVL2.5-26B’s response correctly identifies that no grouping occurs for the same meta-object in the sequence, with the correct answer being ‘D) None’. However, in the CoT response, the model incorrectly selects option C) 2 and 3. While it correctly states that “the second and third images both have dinosaurs wearing sunglasses,” the lack of detailed analysis leads to an inaccurate conclusion. A closer examination reveals a key difference between the images—the dinosaur in image 3 is holding a keyboard instead of a skateboard, which should have prevented the grouping of the two images. This highlights the importance of providing more detailed and unambiguous cues in CoT reasoning.