# Multi-Crit: Benchmarking Multimodal Judges on Pluralistic Criteria-Following

Tianyi Xiong[1*], Yi Ge[1*], Ming Li[1], Zuolong Zhang[1], Pranav Kulkarni[1], Kaishen Wang[1], Qi He[1],
Zeying Zhu[1], Chenxi Liu[1], Ruibo Chen[1], Tong Zheng[1], Yanshuo Chen[1], Xiyao Wang[1],
Renrui Zhang, Wenhu Chen[2], Heng Huang[1]

[1]University of Maryland, College Park, [2]University of Waterloo
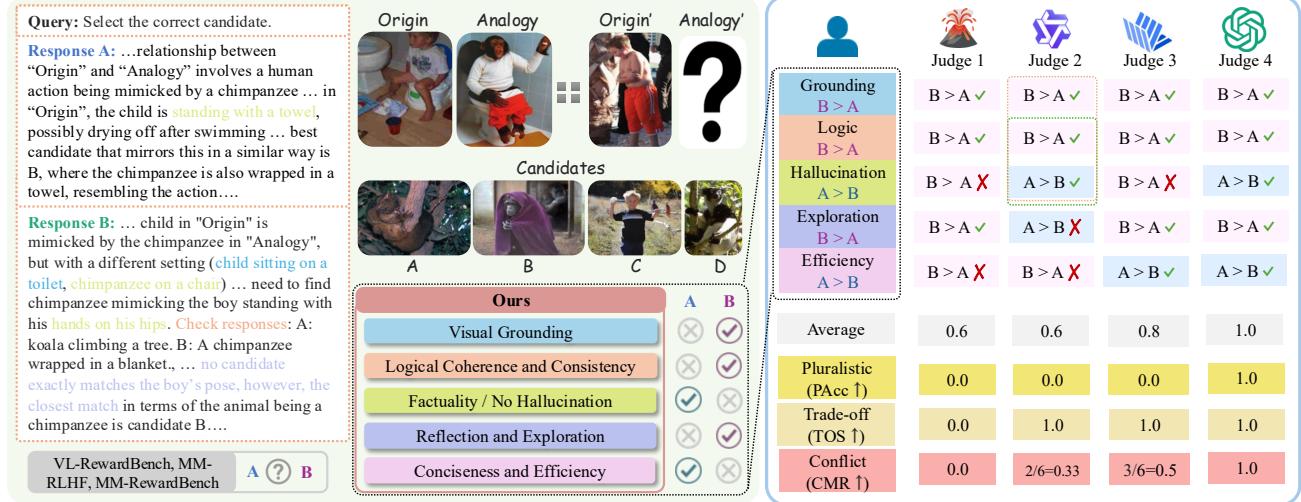
multi-crit.github.io

Figure 1. **Overview of Multi-Crit.** *Left*: Unlike prior works that assign a single overall preference label, Multi-Crit provides pluralistic, multi-criterion human judgments, exposing conflicts between different evaluation criteria within the same sample (e.g., Logic vs. No Hallucination, Reflection vs. Efficiency). *Right*: We introduce three complementary metrics to systematically assess LMM judges on their ability to follow pluralistic evaluation criteria, recognize preference trade-offs, and capture criterion-level conflicts.

## Abstract

*Large multimodal models (LMMs) are increasingly adopted as judges in multimodal evaluation systems due to their strong instruction following and consistency with human preferences. However, their ability to follow diverse, fine-grained evaluation criteria remains underexplored. We develop Multi-Crit, a benchmark for evaluating multimodal judges on their capacity to follow pluralistic criteria and produce reliable criterion-level judgments. Covering both open-ended generation and verifiable reasoning tasks, Multi-Crit is built through a rigorous data curation pipeline that gathers challenging response pairs with multi-criterion human annotations. It further introduces three novel metrics for systematically assessing pluralistic adherence, criterion-switching flexibility, and the ability to recognize criterion-level preference conflicts. Comprehen-*

*sive analysis of 25 LMMs reveals that 1) proprietary models still struggle to maintain consistent adherence to pluralistic criteria–especially in open-ended evaluation; 2) open-source models lag further behind in flexibly following diverse criteria; and 3) critic fine-tuning with holistic judgment signals enhances visual grounding but fails to generalize to pluralistic criterion-level judgment. Additional analyses on reasoning fine-tuning, test-time scaling, and boundary consistency between open-source and proprietary models further probe the limits of current multimodal judges. As a pioneering study, Multi-Crit lays the foundation for building reliable and steerable multimodal AI evaluation.*

## 1. Introduction

The recent advancements of large multimodal models (LMMs) [2, 4, 6, 8, 16, 21, 24] in visual understanding [13, 37, 50, 53] and reasoning [18, 27, 56] have un-

---

*Core contributors.

derscored the growing need for scalable and reliable multi-modal evaluation systems. However, evaluating free-form outputs in open-ended or reasoning-intensive tasks remains inherently challenging, as static or reference-based metrics often fail to generalize and provide limited insight into model behavior. A growing trend has emerged to deploy LMMs themselves as generative judges (*LMM-as-a-Judge*) [46, 57]. Given a multimodal prompt consisting of an image and a user query, model responses, and a set of predefined evaluation criteria, the judge model produces either scalar scores for individual responses or pairwise preference judgments, accompanied by textual justifications for its evaluation. Owing to its scalability and flexibility in defining task-specific criteria, this paradigm has been widely adopted for open-ended evaluation across diverse multimodal benchmarks [13, 18, 21, 26, 34, 45, 49, 52, 56]. A series of works fine-tune open-source models to develop dedicated judge and critic capacities, and have proven effective in providing AI feedback for aligning LMMs through reinforcement fine-tuning [22, 22, 25, 40, 43] or test-time scaling [41, 58].

The reliability of the *LMM-as-a-Judge* paradigm relies on two key components: (1) **consistency with human judgment**, and (2) **flexibility to follow diverse, task-specific evaluation criteria** that capture different aspects of model behavior (e.g., factuality, visual grounding, structural coherence). However, existing multimodal reward and judge benchmarks [5, 23, 55, 60] mainly address the first aspect by providing a single overall preference label for each pair of model responses. This coarse formulation overlooks the evaluation of judge models in producing criterion-specific judgment. Specifically, two responses may exhibit trade-offs across different criteria—for instance, a concise but factually correct answer versus a visually grounded yet slightly hallucinated one. This motivates a key question: *Can LMM-based judges follow diverse evaluation criteria and recognize conflicts among criterion-level judgments?*

To address this limitation, we introduce Multi-Crit, a challenging benchmark designed to comprehensively evaluate multimodal judges on their ability to follow diverse evaluation criteria and produce reliable, criterion-level judgments aligned with human evaluators. As shown in Figure 1, each evaluation sample is annotated with multiple criterion-level human preferences, capturing distinct aspects of response quality and the potential conflicts between criteria. Multi-Crit is built upon two key research questions:

1. **How to construct evaluation data with multi-criterion human preferences that capture diverse aspects of response quality and reveal intra-sample conflicts?** We design a rigorous data curation and preference-annotation pipeline that: (1) collects diverse multimodal prompts spanning both open-ended content generation and verifiable reasoning tasks; (2) selects challenging response pairs from various strong LMMs with fine-grained criterion-level quality differences via multi-step filtering; (3) defines clear evaluation criteria covering complementary aspects of multimodal judgment; and (4) gathers high-quality human annotations ensuring reliability and inter-annotator consistency.

2. **How to assess the performance of LMM-based judges in pluralistic criteria following?** Apart from criterion-level accuracy, we introduce three metrics to systematically assess a judge model's (1) adherence to pluralistic criteria (PAcc), (2) sensitivity to criterion-level trade-offs within each prompt (TOS), and (3) ability to resolve criterion-level preference conflicts (CMR).

Comprehensive evaluation across 25 LMMs exposes key challenges: 1) Strongest proprietary models (e.g., *o4-mini*, *Claude-3.7-Sonnet*) struggle to maintain consistent pluralistic criteria following, achieving 32.78% in open-ended and 53.17% in reasoning judgments; 2) open-source LMMs lag notably behind in both criterion-level accuracy and their ability to recognize trade-offs and preference conflicts; and 3) models fine-tuned on holistic, single-preference critic data primarily improve judgment consistency in visual grounding but generalize poorly to pluralistic, criterion-driven evaluation. Additional analyses further reveals the boundaries of current LMM judges: 1) reasoning fine-tuning fails to enhance reasoning judgment but weakens models' ability to recognize trade-offs; 2) test-time scaling effects are most evident in the strongest model, *o4-mini*, but remain inconsistent across others in pluralistic criteria following; 3) the criterion-level upper bound of proprietary models aligns closely with human inter-annotator consistency, whereas open-source judges fail to exhibit this trend.

Overall, Multi-Crit serves as the first benchmark for assessing multimodal judges' adherence to pluralistic evaluation criteria, built upon a curated dataset, novel metrics, and analyses that uncover key limitations and directions for more reliable, steerable evaluation systems.

## 2. Related Work

**Large Multimodal Model as Judges and Critics.** Large multimodal models (LMMs) are widely adopted for evaluating model-generated responses, providing assessments ranging from scalar [13, 21, 53, 56] to pairwise preference [24, 26, 52]. GPT-4V [57] was first shown to align closely with human judges across diverse multimodal evaluation tasks, motivating efforts to develop open-source counterparts. Subsequent works fine-tuned open-source LMMs into either BT-style reward models with dedicated reward heads [34, 54, 59] or generative judges that produce textual justifications and scores [20, 41, 43, 46, 58, 60] in an autoregressive manner. While our benchmark centers on generative judge models—due to their flexible instruction-following abilities inherent in the autoregressive generation
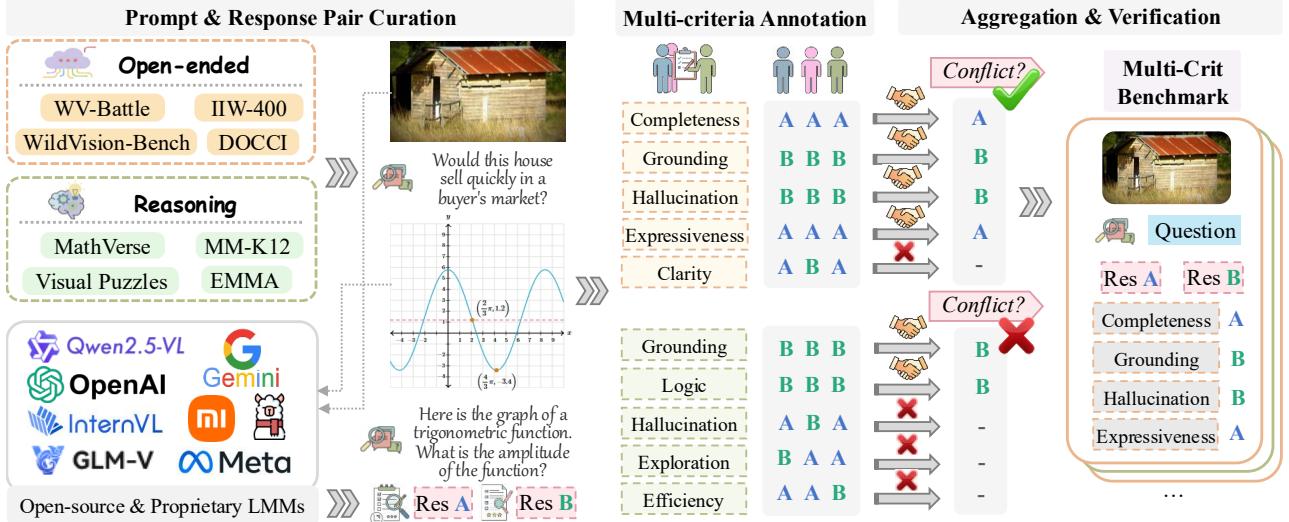
Figure 2. **Data Construction Pipeline.** Multi-Crit is built from diverse prompts across open-ended and reasoning tasks, responses from various LMMs reflecting subtle quality distinctions, and multi-criterion human annotations highlighting preference conflicts across criteria.

process of LMMs—the proposed methodology is broadly applicable across diverse reward modeling paradigms.

**Benchmarking Multimodal Judges.** Benchmarking reward and judge models is crucial for evaluating the reliability of AI feedback, which underpins both automated evaluation and the enhancement of LMM performance through preference alignment. MLLM-as-a-Judge [5] is the first to evaluate LMM as judges on academic tasks; VL-Rewardbench [23], multimodal-rewardbench [48] and MM-RLHF [60] Benchs extend prompts sources to diverse real-world scenarios, reasoning and safety tasks. Following reward and judge benchmarks in text-only LLMs [19, 35, 61], these multimodal benchmarks evaluate judges by their overall preference between paired responses, measured against human or verifier-provided golden labels. MM-Critic further extends evaluation to assess the correctness of generated textual critiques by leveraging GPT-generated scoring rubrics [55]. However, none of these benchmarks evaluate multimodal judges on their ability to follow diverse criteria for each response pair, and fine-grained, criterion-level evaluation remains underexplored.

**Criteria-Following Judgment.** When deploying model-based judges [24, 46, 52, 53, 56], existing benchmarks typically define task-specific criteria, such as helpfulness [26], visual grounding [49], or hallucination detection [34]. While recent work has begun evaluating criterion-following in the textual domain—either by embedding criterion-level differences into responses [32] or summarizing from human justifications [44]—we extend this line of research to the multimodal setting, providing a comprehensive analysis of how judges follow diverse criteria and capture criterion-level preference conflicts.

## 3. Multi-Crit Benchmark

Multi-Crit is a rigorously curated benchmark designed to evaluate the pluralistic criteria-following ability of multimodal judge models. It includes diverse prompts from 8 data sources across open-ended and reasoning-dense domains, cross-model and intra-model response pairs from 11 LMMs of varying capabilities, and fine-grained multi-criteria human annotations within each prompt. The data construction pipeline is illustrated in Figure 2. Additionally, three metrics are introduced in Multi-Crit to evaluate judge models on their ability to flexibly follow diverse evaluation criteria and identify criterion-level preference conflicts.

### 3.1. Task Formulation

Existing multimodal judge benchmarks [5, 23, 48, 60] typically consist of preference pairs $(q, l_a, l_b, y)$, where $q$ is a multimodal prompt containing an image and a textual query, $l_a$ and $l_b$ are candidate responses, and $y \in \{a, b\}$ denotes the golden preference by human annotation or rule-based verification. However, such one-dimensional annotations cannot capture the multifaceted nature of human evaluation—different criteria may lead to conflicting judgments for the same response pair. Multi-Crit therefore extends this formulation to include criterion-specific preferences:

$$(q, l_a, l_b, \{(c_i, y_i)\}_{i=1}^{K_q}),$$

where each $c_i$ denotes the $i$-th evaluation criterion and $y_i \in \{a, b\}$ represents the preferred response under that criterion.

### 3.2. Prompt and Response Pair Curation

MMCrit primarily covers two major scopes of multimodal judgment: (1) *open-ended content generation*, where responses are free-form and traditional fixed metrics are lim-

| Statistic | Number |
|---|---|
| *Open-ended* | |
| - Prompts | 299 |
|   - with conflicting preferences | 206 (68.9%) |
| - Criterion-level judgment | 1,000 |
|   - Conflicting criterion pairs | 501 |
| *Verifiable Reasoning* | |
| - Prompts | 126 |
|   - with conflicting preferences | 109 (86.5%) |
| - Criterion-level judgment | 425 |
|   - Conflicting criterion pairs | 281 |
| Question length quantiles | (5, 10, 27) |
| Response length quantiles | (104, 163.5, 247) |

Table 1. Multi-Crit's key statistics.



Figure 3. Distribution of prompt sources (left) and evaluation criteria (right).

ited; and (2) *reasoning tasks*, where the judge model evaluates the quality of model-generated reasoning processes leading to objectively verifiable answers.

**Prompt Sources.** For open-ended generation, we select images from ImageInWords [11] and DOCCI [28] evaluation set for captioning and incorporate image–question pairs from the WildVision-Bench and -Battle datasets [26], which contain diverse, in-the-wild visual conversation queries. For verifiable reasoning, we draw from reasoning-dense domains, including MathVerse [56] (math), MM-K12 [27] (STEM), EMMA-mini [14] (math, science, coding), and VisualPuzzles [33] (visual analogy and puzzle solving).

**Response Generation and Pair Curation.** To minimize data contamination, we generate candidate responses using 11 high-performing models, covering both proprietary (e.g. GPT-4o [16], Gemini-2.5-Flash [9]) and open-source families (e.g., LLaMA3.2-V [12], Qwen2.5-VL [4], InternVL3 [8]); details in Appendix. We then form two types of response pairs: *cross-model pairs*, generated by randomly selecting two different LMMs, and *intra-model pairs*, produced by the same model via temperature sampling. For the latter, we sample $N = 5$ responses at temperature 1.0 and select the pair with the largest cosine distance between their MiniLM-L6 [38] embeddings to ensure content diversity. This process yields 3,538 response pairs.

**Filtering Challenging Pairs.** We apply a three-stage filtering mechanism to retain challenging yet balanced examples suitable for fine-grained evaluation. *(1) Length normalization:* We exclude pairs with excessive length disparity, mostly retaining those within [0.7, 1.4] to mitigate superficial bias. *(2) Verifiable reasoning filtering:* For reasoning tasks, we use GPT-4o-mini to verify both responses and keep only pairs where both are correct or both incorrect, as correctness alone is trivial and our evaluation targets different criteria of reasoning quality. *(3) Ensemble-based difficulty filtering:* To emphasize nuanced criterion-related differences, we use three strong judge models (GPT-4o, Gemini-2.5-Flash, and Claude-3.7-Sonnet) for initial assessments of overall quality; pairs on which all judges
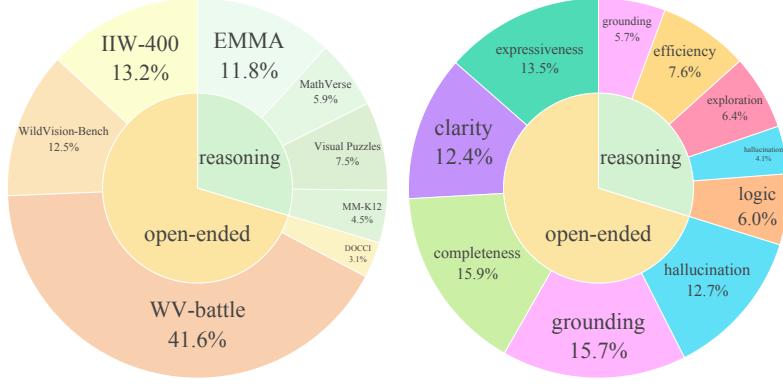
unanimously agree are discarded, while the remaining ones are retained as challenging cases reflecting subtle criterion-level distinctions. A total of 707 response pairs are selected.

### 3.3. Criteria-Specific Preference Annotation

**Criteria Design.** The criteria in Multi-Crit follow a clear set of principles to ensure balanced coverage of multimodal judging capabilities: (1) *Practicality*: reflecting common use cases of multimodal judges; (2) *Specificity*: capturing diverse quality aspects of generated responses while avoiding overlap; and (3) *Generality*: evaluating fundamental dimensions of model behavior rather than prompt-specific content. Following these principles, we survey existing MLLM-as-a-Judge benchmarks, summarize their adopted criteria, and perform multi-round refinement to derive five-criteria for open-ended tasks and another five for reasoning tasks, as briefed in Table 2 and detailed in Appendix.

---

**Criteria for Judging Open-ended Generation**

**Completeness and Coverage:** Address the full scope of the task in the user's query, covering all major elements specified in the prompt as well as relevant visual aspects and contextual cues.
**Visual Grounding and Details:** Reference observable elements in the image such as objects, spatial relationships, colors, or text, and bases its description or analysis on these details.
**Factuality / No Hallucination:** Avoid visual or factual errors, ensuring all details and claims are presented in the image or reasonably supported by the prompt.
**Creativity and Expressiveness:** Demonstrates imagination and originality when appropriate, or precise and knowledgeable articulation for analytical tasks, while remaining contextually appropriate.
**Clarity and Coherence:** Communicates ideas clearly and logically, with fluent language, well-organized structure, and smooth flow of information.

---

**Criteria for Judging Verifiable Reasoning**

**Visual Grounding:** Reference important visual elements—such as objects, layout, or text—and integrates them meaningfully into the reasoning.
**Logic Coherence and Consistency:** Follow a clear, step-by-step logic without contradictions or unjustified leaps, and ensures the answer aligns with reasoning.
**Factuality / No Hallucination:** Ensure accuracy of all claims and support them with the input, avoiding hallucinated visual details or factual errors.
**Reflection and Exploration:** Demonstrate depth of reasoning through reflection and exploration of alternative interpretations, particularly in complex tasks.
**Conciseness and Efficiency:** Remains concise and focused, matching the task complexity while avoiding redundancy or unnecessary over-analysis.

---

Table 2. Summary of the evaluation criteria employed in the Multi-Crit open-ended and reasoning splits, each capturing distinct dimensions of multimodal judgment.

**Rigorous Human Annotation.** We recruit nine Ph.D. students majoring in Computer Science as annotators, all with strong backgrounds in multimodal AI and STEM-related domains. Each receives detailed guidelines with clear illustrations of all evaluation criteria. During annotation, annotators are presented with a multimodal instance containing a prompt and two candidate responses. They exclusively evaluate one criterion at a time, determine which response is better (or declare a tie, limited to under 10%), and provide a brief textual justification with supporting evidence. Annotators focus on response elements most relevant to the prompt and indicative of the given criterion. To ensure consistency, we first establish a seed set of 10 open-ended and 10 reasoning examples. Annotators label this set individually, then participate in group discussions led by the project lead to align interpretations and reach consensus. After calibration, each evaluation sample—comprising the prompt, response pair, and set of criteria—is randomly assigned to three annotators for cross-validation in the main phase. The annotation process involved 289 hours of human effort in total. Inter-annotator agreement is quantified using the average pairwise Cohen's $\kappa$ [3], achieving 0.718 on open-ended generation tasks and 0.805 on verifiable reasoning tasks, indicating substantial and reliable consistency among annotators.

**Preference Aggregation and Final Verification.** First, preferences for each criterion are aggregated, retaining only cases where all annotators agree on one model or where two annotators agree and the third declares a tie. Project leads manually review the textual justifications and discard samples with notable inconsistencies or verbose explanations. Subsequently, criterion-level judgments are aggregated at the prompt–response level. Finally, prompt-level filtering primarily retains pairs exhibiting criterion-level conflicts and balances the number of samples across all criteria.

### 3.4. Data Statistics

Multi-Crit comprises 425 multimodal prompts and 1,425 criterion-level human judgments across open-ended and verifiable reasoning tasks, including 315 prompts with criterion-level conflicts that together constitute 782 conflicting criterion pairs. Key statistics and data distributions are shown in Table 1 and Figure 3, with details in the Appendix.

### 3.5. Measuring Pluralistic Criteria Following

We introduce three complementary metrics to evaluate how well multimodal judge models adhere to pluralistic criteria, identify criterion-level trade-offs, and capture preference conflicts within response pairs.

**Pluralistic Accuracy (PAcc).** Measures whether the judge produces correct judgments across *all* criteria for each evaluation instance. Let $X$ denote the set of evaluation instances, where each $x = (q, l_a, l_b) \in X$ represents a mul-

timodal prompt with its paired candidate responses, and let $C_x$ denote the set of evaluation criteria associated with $x$:

$$\text{PAcc} = \frac{1}{|X|} \sum_{x \in X} \mathbb{I}\left[ \bigwedge_{c \in C_x} \hat{y}_{x,c} = y_{x,c} \right],$$

where $\hat{y}_{x,c}$ denotes the judge's predicted preference and $y_{x,c}$ the ground-truth label for criterion $c$.

**Trade-off Sensitivity (TOS).** Measures whether the judge model recognizes criterion-level trade-offs within each evaluation instance. Let $\mathcal{P}_x = \{ (c_i, c_j) \in C_x \times C_x \mid y_{x,c_i} \neq y_{x,c_j} \}$ denotes the set of criterion pairs in instance $x$ that exhibit ground-truth conflicts.

$$\text{TOS} = \frac{1}{|X^*|} \sum_{x \in X^*} \mathbb{I}\left[ \exists\, (c_i, c_j) \in \mathcal{P}_x : \hat{y}_{x,c_i} \neq \hat{y}_{x,c_j} \right],$$

where $X^* = \{ x \in X \mid \exists\, c_i, c_j \in C_x \text{ such that } y_{x,c_i} \neq y_{x,c_j} \}$ denotes the subset of evaluation instances exhibiting at least one ground-truth conflict among criteria. TOS reflects whether the judge can *perceive* criterion-level trade-offs within an instance, indicating its flexibility rather than absolute accuracy. A higher TOS suggests the model is less criterion-agnostic or overconfident in assigning identical preferences across all criteria.

**Conflict Matching Rate (CMR).** Measures the judge's ability to correctly resolve ground-truth conflicts between criterion-based judgments. A conflict is considered *matched* only if the model predicts both criteria in agreement with the ground truth preference:

$$\text{CMR} = \frac{\sum_{x \in X^*} \sum_{(c_i, c_j) \in \mathcal{P}_x} \mathbb{I}\left[ (\hat{y}_{x,c_i}, \hat{y}_{x,c_j}) = (y_{x,c_i}, y_{x,c_j}) \right]}{\sum_{x \in X^*} |\mathcal{P}_x|}.$$

This provides a *stricter* and more fine-grained metric than TOS, evaluating whether the judge not only detects but also resolves conflicts consistently with humans.

## 4. Experiments

To comprehensively assess multimodal judge boundaries, we evaluate 25 LMMs across three model groups: (1) *proprietary models*, including the GPT series [16, 17, 30], Gemini family [9], and Claude-3.7-Sonnet [31]; (2) *leading open-source LMMs*, including Qwen2.5/3-VL [4, 47], InternVL3/3.5 [39, 62], GLM-4.1V-Thinking [15], LLaMA3.2-V [12], MiniCPM-V-4.5 [51], Eagle2.5 [6], Molmo [10] and LLaVA-OneVision [21]; and (3) *finetuned judge models*, including LLaVA-Critic [46]/-R1 [41], R1-Reward [58] and UnifiedReward [43].

| Model | Prompt-Level | | | Criterion-Level | | | | | Avg |
|---|---|---|---|---|---|---|---|---|---|
| | Pluralisic | Conflict | Tradeoff | Completeness | Grounding | Hallucination | Expressiveness | Clarity | |
| *Open-Source LMMs* | | | | | | | | | |
| Qwen2.5-VL-7B-Instruct | 9.41 | 17.28 | 36.14 | 56.12 | 51.70 | 48.20 | 64.12 | 51.82 | 54.39 |
| LLaVA-OneVision-7B | 11.46 | 12.77 | 37.86 | 53.20 | 48.80 | 43.10 | 55.68 | 52.26 | 50.61 |
| Llama-3.2-11B-Vision-Instruct | 12.71 | 15.37 | 41.26 | 55.75 | 44.64 | 48.07 | 61.14 | 46.59 | 51.24 |
| Molmo-7B | 17.06 | 15.97 | 34.47 | 48.23 | 44.20 | 55.80 | 58.03 | 50.57 | 51.37 |
| Qwen3-VL-8B-Instruct | 18.39 | 18.36 | 34.47 | 56.19 | 47.77 | 59.12 | 55.44 | 52.27 | 54.16 |
| MiniCPM-V-4.5-8B | 18.73 | 14.77 | 27.67 | 53.10 | 46.43 | 53.59 | 65.80 | 50.00 | 53.78 |
| Eagle2.5-8B | 20.40 | 23.75 | 50.97 | 64.16 | 54.91 | 52.49 | 59.59 | 54.55 | 57.14 |
| GLM-4.1V-9B-Thinking | 24.08 | 30.54 | 55.34 | 59.73 | 55.80 | 61.33 | 67.36 | 51.14 | 59.07 |
| InternVL3.5-8B | 25.08 | 32.34 | 62.14 | 61.50 | 59.38 | 56.35 | 69.43 | 58.52 | 61.04 |
| InternVL3-8B-Instruct | 26.09 | 30.34 | 61.17 | 64.60 | 62.95 | 56.35 | 67.88 | 56.25 | 61.61 |
| Qwen2.5-VL-72B-Instruct | 28.43 | 35.53 | 60.68 | 69.47 | 63.39 | 58.56 | <u>70.98</u> | 56.82 | 63.84 |
| MiMo-VL-7B | 29.10 | 39.52 | 65.53 | 65.93 | 62.95 | 64.09 | 70.47 | 53.41 | 63.37 |
| InternVL3-78B-Instruct | 29.10 | 32.53 | 56.31 | <u>73.01</u> | 65.62 | 56.35 | 68.91 | 59.66 | 64.71 |
| InternVL3.5-38B | 30.43 | 33.73 | 64.08 | 71.68 | 65.18 | 62.98 | 63.73 | 61.93 | 65.10 |
| *Finetuned Judge Models* | | | | | | | | | |
| LLaVA-Critic-7B (LlaVA-OV) | 11.70 | 14.17 | 43.2 | 56.93 | 47.32 | 29.82 | 50.28 | 44.08 | 45.69 |
| R1-Reward-7B (Qwen2.5-VL) | 17.73 | 20.36 | 45.63 | 59.29 | 60.71 | 49.72 | 55.44 | 53.98 | 55.83 |
| UnifiedReward-7B(Qwen2.5-VL) | 18.06 | 8.38 | 16.50 | 57.96 | 52.23 | 52.49 | 57.51 | 55.68 | 55.17 |
| LLaVA-Critic-R1-7B(Qwen2.5-VL) | 18.39 | 17.96 | 39.32 | 55.31 | 57.59 | 46.96 | 63.73 | 55.11 | 55.74 |
| *Proprietary LMMs* | | | | | | | | | |
| Gemini-2.5-Flash | 25.42 | 31.34 | 62.14 | 64.6 | 60.71 | 66.30 | 63.73 | 57.39 | 62.55 |
| Gemini-2.5-Pro | 28.76 | 37.92 | **66.50** | 65.93 | 59.82 | 70.17 | 62.18 | 60.23 | 63.67 |
| GPT-5 | 29.77 | 38.52 | 62.62 | 69.91 | 69.64 | **75.69** | 58.55 | **68.75** | 68.51 |
| o3 | 31.10 | 42.71 | 62.62 | 66.37 | **74.55** | <u>72.93</u> | 63.21 | **68.75** | 69.16 |
| GPT-4o | 31.44 | **44.91** | <u>66.02</u> | 68.14 | <u>71.88</u> | 65.75 | **76.17** | <u>65.91</u> | <u>69.57</u> |
| Claude-3.7-Sonnet | <u>31.77</u> | 42.32 | 64.08 | 71.68 | **74.55** | 63.54 | 66.84 | 60.23 | 67.37 |
| o4-mini | **32.78** | <u>43.11</u> | 64.56 | **76.11** | **74.55** | 65.75 | 69.43 | 62.50 | **69.67** |

Table 3. **Results on open-ended split.** Criterion-level judgment performance is summarized by mean accuracy. The best and second best results are shown in **bold** and <u>underlined</u> respectively. Notably, the top performing *o4-mini* only achieves 32.78% in pluralistic accuracy.

**Implementation Details.** During inference, each model was prompted to evaluate a single multimodal sample containing one prompt and a pair of responses, focusing exclusively on one criterion at a time. The full evaluation prompts are provided in the Appendix, with minor adjustments for model-specific formatting. Sampling temperatures follow the original repository settings, or default to 0.6 with a top-$p$ of 0.95 when unspecified.

## 4.1. Main Results

Results are presented in Table 3 for open-ended responses and Table 4 for reasoning responses, revealing the following key observations and analyses:

**Pluralistic judging is highly challenging, especially for open-ended tasks.** The Pluralistic Accuracy (PAcc), which requires correct preferences across all criteria for a single query, is notably low. *o4-mini* and *Claude-3.7-Sonnet* achieve the highest open-ended PAcc at 32.78% and 31.77%, while *o4-mini* and *GPT-5* lead on verifiable reasoning with 53.17% and 45.24%. This highlights the difficulty of pluralistic judging even for state-of-the-art LMMs. As in Figure 4, for top-performing models, results on open-ended
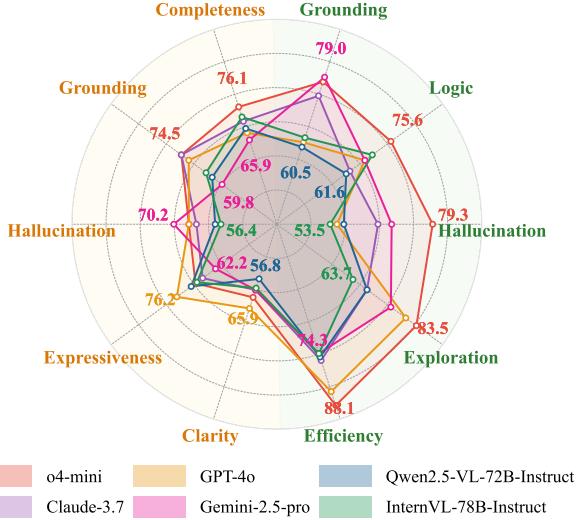


Figure 4. Average performance across each criterion. While the top model differs across criteria, all models show stronger pluralistic alignment in verifiable reasoning than in open-ended tasks.

tasks are lower, reflecting the greater subjectivity and reliance on fine-grained visual perception for evaluating free-form responses (e.g., real-world images or screenshots). In contrast, existing multimodal reasoning tasks—typically in-

| Model | Prompt-Level | | | Criterion-Level | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Pluralisic | Conflict | Tradeoff | Grounding | Logic | Hallucination | Exploration | Efficiency | Avg |
| *Open-Source LMMs* | | | | | | | | | |
| LLaVA-OneVision-7B | 9.52 | 17.08 | 51.38 | 40.74 | 44.19 | 29.31 | 50.55 | 58.72 | 44.70 |
| Llama-3.2-11B-Vision-Instruct | 10.32 | 19.93 | 44.95 | 49.38 | 45.35 | 43.10 | 53.85 | 62.39 | 50.81 |
| Molmo-7B | 11.11 | 13.17 | 34.86 | 45.68 | 40.70 | 43.10 | 43.96 | 61.47 | 46.98 |
| Qwen3-VL-8B-Instruct | 16.67 | 19.22 | 35.78 | 50.62 | 56.98 | 65.52 | 47.25 | 64.22 | 56.92 |
| Qwen2.5-VL-7B-Instruct | 16.67 | 25.62 | 66.06 | 46.91 | 52.33 | 55.17 | 57.14 | 65.14 | 55.34 |
| Eagle2.5-8B | 19.05 | 27.05 | 48.62 | 49.38 | 47.67 | 44.83 | 58.24 | 69.72 | 53.97 |
| GLM-4.1V-9B-Thinking | 23.02 | 27.40 | 55.05 | 54.32 | 60.47 | 65.52 | 50.55 | 66.97 | 59.57 |
| MiniCPM-V-4.5-8B | 25.40 | 38.08 | 75.23 | 56.79 | 66.28 | 56.90 | 65.93 | 66.97 | 62.57 |
| InternVL3-8B-Instruct | 26.98 | 39.50 | 66.06 | 64.20 | 65.12 | 58.62 | 67.03 | 76.15 | 66.22 |
| InternVL3-78B-Instruct | 29.37 | 39.50 | 67.89 | 62.96 | 69.77 | 53.45 | 63.74 | 74.31 | 64.85 |
| Qwen2.5-VL-72B-Instruct | 32.54 | 45.91 | 77.06 | 60.49 | 61.63 | 56.90 | 68.13 | 75.23 | 64.48 |
| InternVL3.5-8B | 32.54 | 39.15 | 69.72 | 65.43 | 63.95 | 56.90 | 71.43 | 71.56 | 65.85 |
| MiMo-VL-7B | 37.30 | 41.99 | 71.56 | 59.26 | 59.30 | 70.69 | 67.03 | 75.23 | 66.30 |
| InternVL3.5-38B | 37.30 | 47.69 | 75.23 | 71.60 | <u>70.93</u> | 55.17 | 78.02 | 73.39 | 69.82 |
| *Finetuned Judges* | | | | | | | | | |
| UnifiedReward-7B (Qwen2.5-VL) | 11.90 | 13.17 | 22.02 | 64.2.0 | <u>70.93</u> | 50.00 | 58.24 | 54.13 | 59.50 |
| LLaVA-Critic-7B (LLaVA-OV) | 13.49 | 22.42 | 50.46 | 43.21 | 44.19 | 43.10 | 45.05 | 60.55 | 47.22 |
| LLaVA-Critic-R1-7B (Qwen2.5-VL) | 16.67 | 16.37 | 49.54 | 53.09 | 50.00 | 51.72 | 56.04 | 55.05 | 53.18 |
| R1-Reward (Qwen2.5-VL) | 19.05 | 24.56 | 62.39 | 58.02 | 54.65 | 46.55 | 52.75 | 60.55 | 54.50 |
| *Proprietary Judge Models* | | | | | | | | | |
| Claude-3.7-Sonnet | 36.51 | 51.96 | <u>83.49</u> | 74.07 | 62.79 | 65.52 | 68.13 | 76.15 | 69.33 |
| Gemini-2.5-Flash | 37.30 | 53.05 | 73.39 | 67.50 | 63.95 | 64.91 | 78.02 | 77.98 | 70.47 |
| Gemini-2.5-Pro | 41.27 | 52.33 | 75.93 | **79.01** | 67.44 | 68.97 | 75.56 | 74.31 | 73.06 |
| GPT-4o | 41.27 | 55.16 | **84.40** | 61.73 | 67.44 | 55.17 | 80.22 | <u>84.40</u> | 69.79 |
| o3 | 44.44 | <u>62.28</u> | 82.57 | 67.90 | <u>70.93</u> | **84.21** | **84.62** | 81.65 | <u>77.86</u> |
| GPT-5 | <u>45.24</u> | 56.58 | 78.90 | 74.07 | <u>70.93</u> | <u>82.76</u> | 81.32 | 77.98 | 77.41 |
| o4-mini | **53.17** | **65.84** | <u>83.49</u> | <u>77.78</u> | **75.58** | 79.31 | <u>83.52</u> | **88.07** | **80.85** |

Table 4. **Results on verifiable reasoning.** Best and second-best results are in **bold** and <u>underlined</u>.

volving multimodal math, STEM, puzzle, and coding problems—are more objective, with clearer verification and explicit visual grounding. These findings highlight the inherent difficulty of pluralistic judgment and underscore the necessity of the proposed Multi-Crit benchmark.

**Model performance varies across criteria and metrics.**
For judgment in verifiable reasoning, *o4-mini* performs best on Logic (75.58%) and Efficiency (88.07%), but is surpassed by *o3* in Hallucination (84.21%) and Exploration (84.62%), and by *Gemini-2.5-Pro* in Grounding (79.01%). At the prompt level, *GPT-4o* shows the highest sensitivity to trade-offs in prompts with criterion conflicts, achieving the best Trade-off Sensitivity of 84.40%. However, its ability to resolve conflicting criterion-level preference pairs is notably weaker than that of the top-performing *o4-mini* (55.16% vs. 65.84% in CMR). A similar trend is observed in open-ended tasks: *GPT-4o-* excels in *Expressiveness* (76.17%) but lags in *Completeness* (68.14%) and *Hallucination* (65.75%). In contrast, *o3* performs strongly on *Hallucination* (72.93%) but ranks low in *Expressiveness* (63.21%). These results highlight the multifaceted nature of pluralistic judgment—no model excels across all aspects, each showing distinct strengths and blind spots.

**Open-source LMMs fall short in pluralistic criterion-following, especially in capturing criterion-level preference conflicts.** As shown in Tables 3 and 4, proprietary LMMs (e.g., *o4-mini*, *GPT-5*, *o3*) maintain a clear lead over open-source counterparts (e.g., *InternVL3.5-38B*, *Qwen2.5-VL-72B-Instruct*) across all prompt-level metrics (PAcc, CMR, and TOS) and criterion-level accuracies. Specifically, the conflict matching rate drops by about 9.4 points in open-ended tasks (43.1 → 33.7) and 18.1 points in reasoning tasks (65.8 → 47.7), far exceeding the 4- and 11-point declines observed in criterion-level accuracy. This indicates that open-source models struggle to maintain consistent pluralistic evaluation across conflicting criteria, highlighting their limited ability to integrate multi-faceted criterion-level signals into reliable judgments.

**Finetuned judge models improve evaluation in visual grounding but are limited in handling pluralistic criteria.** A striking observation is that models specifically finetuned as judges or reward models (e.g., *LLaVA-Critic* [46], *R1-Reward* [58], *LLaVA-Critic-R1* [41], *UnifiedReward-Think* [43]) do not exhibit consistent performance gains over their base models, unlike prior benchmarks [23, 35]. For Qwen-based fine-tuned judges, *R1-Reward* is the only model that outperforms the base *Qwen2.5-VL-7B-Instruct*

in *PAcc* across both the open-ended split (9.41→17.73) and the reasoning split (16.67→19.05)—yet its improvement at recognizing trade-off (36.14→45.63, 66.06→62.39) and conflict capturing (17.28→20.36, 25.62→24.56) remains inconsistent. Other fine-tuned judges perform similarly or worse than similar-sized models without critic tuning, remaining far behind proprietary counterparts.

Another notable trend is that all Qwen-based fine-tuned judges show consistent improvements in evaluating visual grounding across both splits— *R1-Reward* (51.7→60.7, 46.9→58.0), *LLaVA-Critic-R1* (51.7→57.6, 46.9→53.1), and *UnifiedReward* (51.7→52.2, 46.9→64.2). For the LLaVA-based judge, *LLaVA-Critic-7B (LLaVA-OV)* improves on reasoning (40.7→43.2) but a marginal drop on open-ended evaluation (48.8→47.3).

Overall, current critic fine-tuning pipelines—driven by holistic preference signals and limited critic prompt templates—mainly enhance judgment capacity in visual–textual grounding and alignment, but fail to generalize to diverse, fine-grained, and often conflicting evaluation criteria, thereby underscoring the need for pluralistic, criterion-aware training as introduced in Multi-Crit.

### 4.2. Additional Discussions

To further investigate the capabilities and limitations of LMM judges in pluralistic criteria-following, we conduct three additional analyses as follows:

**1. Does RLVR on multimodal reasoning improve judgment in reasoning traces?** We evaluate three GRPO-finetuned models [7, 27, 42] on general reasoning tasks (e.g., math and charts), all showing clear reasoning gains. However, these models, trained on holistic accuracy rewards, show reduced sensitivity to trade-offs and conflict matching in reasoning judgments (Figure 5), indicating weaker recognition of criterion-level preference conflicts.

Figure 5. Results of RL-tuned reasoning models on the Multi-Crit reasoning split, all based on Qwen2.5-VL-7B.

**2. Does test-time scaling hold?** We run each judge model $K$ times, take a majority vote for each criterion sample as the final judgment, and then compute the pluralistic accuracy. As shown in Figure 6, the strongest model, *o4-mini*, exhibits the most robust test-time scaling behavior for both open-ended (32.78→37.12) and reasoning (53.17→57.94) judgments, with *GPT-4o* showing similar gains. Other models show inconsistent trends and large variance, indicat-

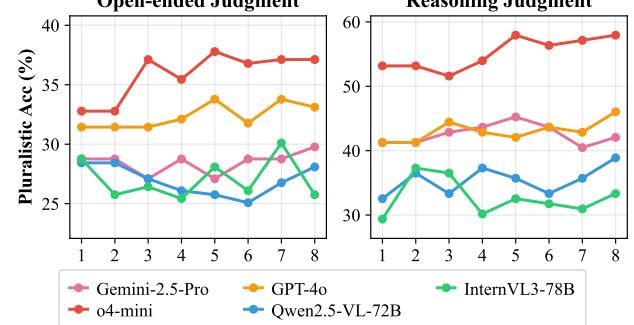ing that test-time scaling offers only limited and model-dependent benefits in multi-criteria judgment.

Figure 6. Test-time scaling behavior by pluralistic accuracy.

**3. What are the current boundaries of criteria following for open-source and proprietary models?** To probe the upper limits of criterion-level consistency, we correlate each model group's strongest criterion accuracies with human inter-annotator agreement (Cohen's $\kappa$), excluding *reasoning-efficiency*, which is largely influenced by response length. As in Figure 7, proprietary models exhibit a strong and significant correlation (r=0.73, p=0.024), indicating that their upper-bound judgment patterns align closely with human agreement. In contrast, open-source models show a weaker, more inconsistent relationship (r=0.36, p=0.344), reflecting limited capacity to internalize human-consistent evaluation criteria. These findings suggest that while open-source models would benefit from scaling high-quality human annotations to improve pluralistic criteria-following, proprietary systems might face the next challenge of surpassing human-level evaluation alignment.
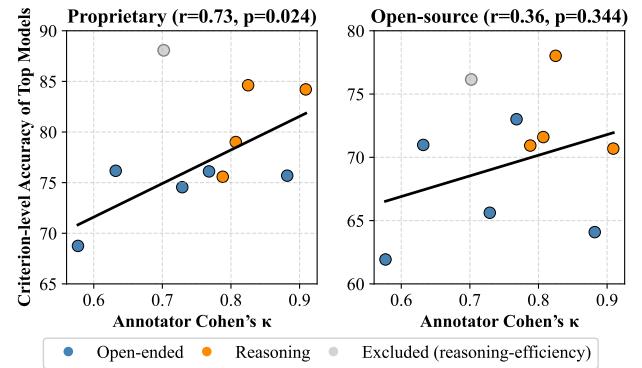
Figure 7. Correlation between criterion-level top judge accuracy and inter-annotator agreement. Each point denotes one criterion.

## 5. Conclusion

We introduced Multi-Crit, a comprehensive benchmark for evaluating multimodal judges under pluralistic, criterion-level settings. By providing fine-grained multi-criterion human annotations and complementary pluralistic metrics,

Multi-Crit enables systematic analysis of models' adherence to diverse criteria, trade-off recognition, and conflict resolution. Comprehensive experiments reveal key limitations of existing *LMM-as-a-Judge* systems: proprietary models struggle with consistent pluralistic adherence, particularly in open-ended evaluation; open-source models lag further behind in their criteria-following capacities; and critic fine-tuning, while improving visual grounding and reference alignment, remains narrowly effective and fails to generalize to fine-grained or conflicting evaluation dimensions. Additional analyses highlight diminished trade-off recognition in reasoning fine-tuning, inconsistent test-time scaling behaviors, and the human-aligned upper bounds of proprietary models in criterion-level consistency. In summary, Multi-Crit serves as a challenging suite to explore the boundaries of current AI-based multimodal judge systems, paving the way toward more steerable, reliable, and even superhuman AI feedback.

## References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 12

[2] Xiang An, Yin Xie, Kaicheng Yang, Wenkang Zhang, Xiuwei Zhao, Zheng Cheng, Yirui Wang, Songcen Xu, Changrui Chen, Chunsheng Wu, et al. Llava-onevision-1.5: Fully open framework for democratized multimodal training. *arXiv preprint arXiv:2509.23661*, 2025. 1

[3] Ron Artstein and Massimo Poesio. Inter-coder agreement for computational linguistics. *Computational linguistics*, 34(4): 555–596, 2008. 5

[4] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 1, 4, 5, 12

[5] Dongping Chen, Ruoxi Chen, Shilin Zhang, Yaochen Wang, Yinuo Liu, Huichi Zhou, Qihui Zhang, Yao Wan, Pan Zhou, and Lichao Sun. Mllm-as-a-judge: Assessing multimodal llm-as-a-judge with vision-language benchmark. In *Forty-first International Conference on Machine Learning*, 2024. 2, 3, 12

[6] Guo Chen, Zhiqi Li, Shihao Wang, Jindong Jiang, Yicheng Liu, Lidong Lu, De-An Huang, Wonmin Byeon, Matthieu Le, Tuomas Rintamaki, et al. Eagle 2.5: Boosting long-context post-training for frontier vision-language models. *arXiv preprint arXiv:2504.15271*, 2025. 1, 5

[7] Hardy Chen, Haoqin Tu, Fali Wang, Hui Liu, Xianfeng Tang, Xinya Du, Yuyin Zhou, and Cihang Xie. Sft or rl? an early investigation into training r1-like reasoning large vision-language models, 2025. 8

[8] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 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, 2024. 1, 4

[9] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 4, 5, 12

[10] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. Molmo and pixmo: Open weights and open data for state-of-the-art vision-language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 91–104, 2025. 5

[11] Roopal Garg, Andrea Burns, Burcu Karagol-Ayan, Yonatan Bitton, Ceslee Montgomery, Yasumasa Onoe, Andrew Bunner, Ranjay Krishna, Jason Michael Baldridge, and Radu Soricut. Imageinwords: Unlocking hyper-detailed image descriptions. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 93–127, 2024. 4

[12] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 4, 5, 12

[13] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14375–14385, 2024. 1, 2

[14] Yunzhuo Hao, Jiawei Gu, Huichen Will Wang, Linjie Li, Zhengyuan Yang, Lijuan Wang, and Yu Cheng. Can mllms reason in multimodality? emma: An enhanced multimodal reasoning benchmark. *arXiv preprint arXiv:2501.05444*, 2025. 4

[15] Wenyi Hong, Wenmeng Yu, Xiaotao Gu, Guo Wang, Guobing Gan, Haomiao Tang, Jiale Cheng, Ji Qi, Junhui Ji, Lihang Pan, et al. Glm-4.1 v-thinking: Towards versatile multimodal reasoning with scalable reinforcement learning. *arXiv e-prints*, pages arXiv–2507, 2025. 5, 12

[16] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 1, 4, 5

[17] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024. 5

[18] Dongzhi Jiang, Renrui Zhang, Ziyu Guo, Yanwei Li, Yu Qi, Xinyan Chen, Liuhui Wang, Jianhan Jin, Claire Guo, Shen Yan, et al. Mme-cot: Benchmarking chain-of-thought in large multimodal models for reasoning quality, robustness, and efficiency. *arXiv preprint arXiv:2502.09621*, 2025. 1, 2

[19] Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward models for language modeling. *arXiv preprint arXiv:2403.13787*, 2024. 3

[20] Seongyun Lee, Seungone Kim, Sue Park, Geewook Kim, and Minjoon Seo. Prometheus-vision: Vision-language model as a judge for fine-grained evaluation. In *Findings of the association for computational linguistics ACL 2024*, pages 11286–11315, 2024. 2

[21] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024. 1, 2, 5

[22] Lei Li, Zhihui Xie, Mukai Li, Shunian Chen, Peiyi Wang, Liang Chen, Yazheng Yang, Benyou Wang, Lingpeng Kong, and Qi Liu. Vlfeedback: A large-scale ai feedback dataset for large vision-language models alignment. *arXiv preprint arXiv:2410.09421*, 2024. 2

[23] Lei Li, Yuancheng Wei, Zhihui Xie, Xuqing Yang, Yifan Song, Peiyi Wang, Chenxin An, Tianyu Liu, Sujian Li, Bill Yuchen Lin, et al. Vl-rewardbench: A challenging benchmark for vision-language generative reward models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24657–24668, 2025. 2, 3, 7, 12

[24] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023. 1, 2, 3

[25] Jie Liu, Gongye Liu, Jiajun Liang, Ziyang Yuan, Xiaokun Liu, Mingwu Zheng, Xiele Wu, Qiulin Wang, Menghan Xia, Xintao Wang, et al. Improving video generation with human feedback. *arXiv preprint arXiv:2501.13918*, 2025. 2

[26] Yujie Lu, Dongfu Jiang, Wenhu Chen, William Yang Wang, Yejin Choi, and Bill Yuchen Lin. Wildvision: Evaluating vision-language models in the wild with human preferences. *Advances in Neural Information Processing Systems*, 37:48224–48255, 2024. 2, 3, 4, 12

[27] Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Botian Shi, Wenhai Wang, Junjun He, Kaipeng Zhang, et al. Mm-eureka: Exploring visual aha moment with rule-based large-scale reinforcement learning. *arXiv preprint arXiv:2503.07365*, 2025. 1, 4, 8

[28] Yasumasa Onoe, Sunayana Rane, Zachary Berger, Yonatan Bitton, Jaemin Cho, Roopal Garg, Alexander Ku, Zarana Parekh, Jordi Pont-Tuset, Garrett Tanzer, et al. Docci: Descriptions of connected and contrasting images. In *European Conference on Computer Vision*, pages 291–309. Springer, 2024. 4

[29] OpenAI. ntroducing gpt-4.1 in the api. `https://openai.com/index/gpt-4-1/`, 2024. Accessed: 2025-11-03. 12

[30] OpenAI. Openai o3 and o4-mini system card. Webpage, 2025. 5

[31] Anthropic PBC. Claude 3.7 sonnet and claude code. `https://www.anthropic.com/news/claude-3-7-sonnet`, 2025. Model release; hybrid reasoning model supporting "extended thinking mode". 5

[32] Silviu Pitis, Ziang Xiao, Nicolas Le Roux, and Alessandro Sordoni. Improving context-aware preference modeling for language models. *Advances in Neural Information Processing Systems*, 37:70793–70827, 2024. 3

[33] Yueqi Song, Tianyue Ou, Yibo Kong, Zecheng Li, Graham Neubig, and Xiang Yue. Visualpuzzles: Decoupling multimodal reasoning evaluation from domain knowledge. *arXiv preprint arXiv:2504.10342*, 2025. 4

[34] Zhiqing Sun, Sheng Shen, Shengcao Cao, Haotian Liu, Chunyuan Li, Yikang Shen, Chuang Gan, Liang-Yan Gui, Yu-Xiong Wang, Yiming Yang, et al. Aligning large multimodal models with factually augmented rlhf. *arXiv preprint arXiv:2309.14525*, 2023. 2, 3

[35] Sijun Tan, Siyuan Zhuang, Kyle Montgomery, William Y Tang, Alejandro Cuadron, Chenguang Wang, Raluca Ada Popa, and Ion Stoica. Judgebench: A benchmark for evaluating llm-based judges. *arXiv preprint arXiv:2410.12784*, 2024. 3, 7

[36] Kimi Team, Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chenzhuang Du, Chu Wei, et al. Kimi-vl technical report. *arXiv preprint arXiv:2504.07491*, 2025. 12

[37] Peter Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Adithya Jairam Vedagiri IYER, Sai Charitha Akula, Shusheng Yang, Jihan Yang, Manoj Middepogu, Ziteng Wang, et al. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *Advances in Neural Information Processing Systems*, 37:87310–87356, 2024. 1

[38] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in neural information processing systems*, 33:5776–5788, 2020. 4

[39] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, et al. Internvl3. 5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025. 5, 24

[40] Xiyao Wang, Jiuhai Chen, Zhaoyang Wang, Yuhang Zhou, Yiyang Zhou, Huaxiu Yao, Tianyi Zhou, Tom Goldstein, Parminder Bhatia, Taha Kass-Hout, et al. Enhancing visual-language modality alignment in large vision language models via self-improvement. In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 268–282, 2025. 2

[41] Xiyao Wang, Chunyuan Li, Jianwei Yang, Kai Zhang, Bo Liu, Tianyi Xiong, and Furong Huang. Llava-critic-r1: Your critic model is secretly a strong policy model. *arXiv preprint arXiv:2509.00676*, 2025. 2, 5, 7, 12

[42] Xiyao Wang, Zhengyuan Yang, Chao Feng, Hongjin Lu, Linjie Li, Chung-Ching Lin, Kevin Lin, Furong Huang, and Lijuan Wang. Sota with less: Mcts-guided sample selection for data-efficient visual reasoning self-improvement. *arXiv preprint arXiv:2504.07934*, 2025. 8

[43] Yibin Wang, Yuhang Zang, Hao Li, Cheng Jin, and Jiaqi Wang. Unified reward model for multimodal understanding and generation. *arXiv preprint arXiv:2503.05236*, 2025. 2, 5, 7, 12

[44] Zhilin Wang, Jiaqi Zeng, Olivier Delalleau, Ellie Evans, Daniel Egert, Hoo-Chang Shin, Felipe Soares, Yi Dong, and Oleksii Kuchaiev. Rlbff: Binary flexible feedback to bridge between human feedback & verifiable rewards. *arXiv preprint arXiv:2509.21319*, 2025. 3

[45] Shaoyuan Xie, Lingdong Kong, Yuhao Dong, Chonghao Sima, Wenwei Zhang, Qi Alfred Chen, Ziwei Liu, and Liang Pan. Are vlms ready for autonomous driving? an empirical study from the reliability, data, and metric perspectives. *arXiv preprint arXiv:2501.04003*, 2025. 2

[46] Tianyi Xiong, Xiyao Wang, Dong Guo, Qinghao Ye, Haoqi Fan, Quanquan Gu, Heng Huang, and Chunyuan Li. Llava-critic: Learning to evaluate multimodal models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13618–13628, 2025. 2, 3, 5, 7, 12

[47] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025. 5, 24

[48] Michihiro Yasunaga, Luke Zettlemoyer, and Marjan Ghazvininejad. Multimodal rewardbench: Holistic evaluation of reward models for vision language models. *arXiv preprint arXiv:2502.14191*, 2025. 3, 12

[49] Qinghao Ye, Xianhan Zeng, Fu Li, Chunyuan Li, and Haoqi Fan. Painting with words: Elevating detailed image captioning with benchmark and alignment learning. *arXiv preprint arXiv:2503.07906*, 2025. 2, 3

[50] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *European conference on computer vision*, pages 69–85. Springer, 2016. 1

[51] Tianyu Yu, Zefan Wang, Chongyi Wang, Fuwei Huang, Wenshuo Ma, Zhihui He, Tianchi Cai, Weize Chen, Yuxiang Huang, Yuanqian Zhao, et al. Minicpm-v 4.5: Cooking efficient mllms via architecture, data, and training recipe. *arXiv preprint arXiv:2509.18154*, 2025. 5

[52] Tianyu Yu, Haoye Zhang, Qiming Li, Qixin Xu, Yuan Yao, Da Chen, Xiaoman Lu, Ganqu Cui, Yunkai Dang, Taiwen He, et al. Rlaif-v: Open-source ai feedback leads to super gpt-4v trustworthiness. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 19985–19995, 2025. 2, 3

[53] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023. 1, 2, 3

[54] Yuhang Zang, Xiaoyi Dong, Pan Zhang, Yuhang Cao, Ziyu Liu, Shengyuan Ding, Shenxi Wu, Yubo Ma, Haodong Duan, Wenwei Zhang, et al. Internlm-xcomposer2. 5-reward: A simple yet effective multi-modal reward model. *arXiv preprint arXiv:2501.12368*, 2025. 2

[55] Gailun Zeng, Ziyang Luo, Hongzhan Lin, Yuchen Tian, Kaixin Li, Ziyang Gong, Jianxiong Guo, and Jing Ma. MM-CRITIC: A holistic evaluation of large multimodal models as multimodal critique. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 13603–13630, Suzhou, China, 2025. Association for Computational Linguistics. 2, 3

[56] Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu Qiao, et al. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? In *European Conference on Computer Vision*, pages 169–186. Springer, 2024. 1, 2, 3, 4

[57] Xinlu Zhang, Yujie Lu, Weizhi Wang, An Yan, Jun Yan, Lianke Qin, Heng Wang, Xifeng Yan, William Yang Wang, and Linda Ruth Petzold. Gpt-4v (ision) as a generalist evaluator for vision-language tasks. *arXiv preprint arXiv:2311.01361*, 2023. 2

[58] Yi-Fan Zhang, Xingyu Lu, Xiao Hu, Chaoyou Fu, Bin Wen, Tianke Zhang, Changyi Liu, Kaiyu Jiang, Kaibing Chen, Kaiyu Tang, et al. R1-reward: Training multimodal reward model through stable reinforcement learning. *arXiv preprint arXiv:2505.02835*, 2025. 2, 5, 7, 12

[59] Yi-Fan Zhang, Haihua Yang, Huanyu Zhang, Yang Shi, Zezhou Chen, Haochen Tian, Chaoyou Fu, Haotian Wang, Kai Wu, Bo Cui, et al. Basereward: A strong baseline for multimodal reward model. *arXiv preprint arXiv:2509.16127*, 2025. 2

[60] Yi-Fan Zhang, Tao Yu, Haochen Tian, Chaoyou Fu, Peiyan Li, Jianshu Zeng, Wulin Xie, Yang Shi, Huanyu Zhang, Junkang Wu, et al. Mm-rlhf: The next step forward in multimodal llm alignment. *arXiv preprint arXiv:2502.10391*, 2025. 2, 3, 12

[61] Enyu Zhou, Guodong Zheng, Binghai Wang, Zhiheng Xi, Shihan Dou, Rong Bao, Wei Shen, Limao Xiong, Jessica Fan, Yurong Mou, et al. Rmb: Comprehensively benchmarking reward models in llm alignment. *arXiv preprint arXiv:2410.09893*, 2024. 3

[62] Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, et al. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*, 2025. 5, 12

## A. Overview of the Appendix

- Section B illustrates the differences between Multi-Crit and prior multimodal judge benchmarks.
- Section C provides additional benchmark details, including the response models, evaluation criteria, judge prompt template, and supplementary data statistics.
- Section D presents qualitative examples of multimodal prompts and complete evaluation instances.
- Section E reports additional results and analysis.

## B. Comparison Against Prior Benchmarks

As shown in Table 5, earlier benchmarks assign only a single overall preference label for each pair of model responses, whereas Multi-Crit provides multiple criterion-level human annotations. This makes Multi-Crit the first benchmark capable of assessing whether LMM judges can follow pluralistic evaluation criteria and deliver reliable criterion-level judgments. In addition, the average length of candidate model responses in Multi-Crit is significantly longer (164 vs. 99 words), further highlighting the increased difficulty and overall challenge of the benchmark.

## C. Benchmark Construction and Statistics

### C.1. Additional Benchmark Details

**Models for Generating Candidate Responses:** Gemini-2.5-Pro [9], Gemini-2.5-Flash [9], GPT-4.1 [29], GPT-4o [1], Qwen2.5-VL-7B-Instruct [4], Qwen2.5-VL-72B-Instruct [4], LLaMA-3.2-11B-Vision-Instruct [12], InternVL3-8B-Instruct [62], InternVL3-38B-Instruct [62], MiMo-VL-7B-RL [36], GLM-4.1V-9B-Thinking [15].

Specifically, for prompts from WildVision-Battle [26], if the two original responses receive a tie in their human preference, we keep the original pair; otherwise, we generate new responses using the models listed above.

**Detailed Evaluation Criteria.** Table 7 presents the detailed criteria for judging open-ended generation, and Table 8 lists those for judging verifiable reasoning. To better guide both human annotators and judge LMMs, we organize the description of each criterion into three components:

1. **Definition**: A general descriptions of the criterion and the key aspects a response should demonstrate.
2. **Positive Indicators**: Behaviors and qualities that are encouraged and should be preferred in responses.
3. **Negative Indicators**: Errors, undesirable behaviors, or deficiencies that should be penalized under this criterion.

**Evaluation Prompt Template.** We provide the evaluation prompt used for judge-model inference in Table 9. When evaluating LMM judges on the Multi-Crit benchmark, each criterion instance—consisting of the question, the pair of responses, and the target criterion—is filled into this template and sent to the model during inference. For

| Benchmark | Resp. Quantile | Multi-Criterion Anno. |
|---|---|---|
| MLLM-as-a-Judge [5] | (54, 89, 153) | ✗ |
| VL-RewardBench [23] | (48, 99, 136) | ✗ |
| MM-RLHF [60] | (16, 49, 144) | ✗ |
| Multimodal-RewardBench [48] | (1, 72, 138) | ✗ |
| Multi-Crit (ours) | (104, 164, 247) | ✓ |

Table 5. Comparison with existing benchmarks.

| Split | #Prompts | Cross-model pairs | Intra-model pairs |
|---|---|---|---|
| Open-ended | 299 | 168 (56.2%) | 131 (43.8%) |
| Reasoning | 126 | 76 (60.3%) | 50 (39.7%) |
| Total | 425 | 244 (57.4%) | 181 (42.6%) |

Table 6. Composition of response pairs.

finetuned judges [41, 43, 46, 58], we adapt this template to match the prompt format used in their original repositories, avoiding significant out-of-distribution judgment behaviors.

### C.2. Additional Data Statistics

**Response Pair Composition.** As shown in Table 6, for both splits, roughly 40% of response pairs originate from the same model via random sampling, while the remaining pairs are generated by different models. This balanced mix ensures that Multi-Crit evaluates whether LMM judges can capture response-quality differences both across models (reflecting model and training diversity) and within the same model (reflecting sampling diversity).

**Cross-Criterion Correlation and Conflict Frequency.** We report additional statistics on criterion-level human preferences to illustrate how Multi-Crit captures cross-criterion correlations and conflicts. Figure 8 shows how human preferences correlate across criteria on prompts that exhibit criterion-level conflicts. Figure 9 reports how often two criteria yield opposite human preferences on the same response pair. Multi-Crit contains preference conflicts across all criterion pairs in both splits, reflecting both the nuanced differences within the response pairs and the inherent diversity and tension across evaluation dimensions. Certain criterion pairs display intuitive trade-offs—such as *completeness vs. no-hallucination* and *no-hallucination vs. expressiveness* in the open-ended split, and *exploration vs. efficiency* in the reasoning split—and these pairs also exhibit the highest frequencies of human-annotated conflicts.

| Evaluation Criterion | Descrption |
| --- | --- |
| **Completeness and Coverage** | The response should provide a thorough and well-developed answer that fully addresses the intent of the prompt. It must address the complete scope of the task, incorporating all major elements specified in the prompt, as well as relevant visual aspects and broader contextual cues. When appropriate, drawing on relevant external knowledge is encouraged to enrich the explanation.<br>- Reward: The response addresses all key parts of the prompt and image, showing depth and effort in the description or explanation.<br>- Penalize: The response is underdeveloped, fails to meet one or more specific requirements in the prompt, or omits important visual elements or interpretive points. |
| **Visual Grounding and Details** | The response should demonstrate a clear and meaningful connection to the visual input. It should refer to observable elements in the image—such as objects, spatial relationships, colors, or text—and build its description or interpretation based on those elements.<br>- Reward: The response explicitly references relevant visual details that are clearly visible in the image.<br>- Penalize: The response fails to connect meaningfully to the image, or uses vague, generic language that lacks specific visual grounding. |
| **Factuality / No Hallucination** | The response should avoid introducing any visual details, objects, relationships, or factual claims that are not present in the image or reasonably suggested by the prompt. This includes both visual hallucinations (e.g., describing elements not visible in the image) and factual inaccuracies in general knowledge.<br>- Reward: The response stays grounded in the image and prompt, without inventing visual elements or making unsupported factual claims.<br>- Penalize: The response introduces hallucinated visual content or inaccurate factual statements that are unsupported or misleading. |
| **Creativity and Expressiveness** | The response should show originality or stylistic flair for open-ended tasks, and knowledge-informed articulation with precision and depth for analytical tasks. All responses must remain contextually appropriate and grounded in the visual input, while enhancing richness, nuance, and overall engagement.<br>- Reward: The response uses vivid language, unique phrasing, or inventive associations that enrich the interpretation, or it demonstrates professional articulation through deep and knowledge-grounded analysis.<br>- Penalize: The response is overly literal, flat, or dull, lacking originality, variation in expression, or in analytical contexts, fails to demonstrate professional depth or expertise. |
| **Clarity and Coherence** | The response should communicate ideas clearly and logically, with coherent structure and fluent language. This involves not only grammatical correctness, but also effective organization of information, smooth transitions, and consistent flow of ideas.<br>- Reward: The response is clearly written, logically structured, and easy to follow. A brief summary at the beginning may further improve clarity.<br>- Penalize: The response is difficult to follow due to unclear structure, disorganized reasoning, poor transitions, or awkward and repetitive phrasing. |

Table 7. Five evaluation criteria for judging *open-ended content generation* tasks in Multi-Crit Benchmark.

| Evaluation Criterion | Descrption |
| --- | --- |
| **Visual Grounding** | The response should be explicitly grounded in the visual input. It must refer to salient visual elements—such as specific objects, spatial arrangements, colors, or visible text—and incorporate them meaningfully into the reasoning process. Visual references should be accurate and relevant to the task.<br>- Reward: The response clearly references important visual features and integrates them into the reasoning in a precise and relevant manner.<br>- Penalize: The response fails to reference relevant visual elements, or uses generic or weakly connected visual details that do not meaningfully support the reasoning. |
| **Logic Coherence and Consistency** | The reasoning should follow a logically sound and step-by-step progression, with each step building upon the previous one. The reasoning should be internally consistent, with no contradictions, missing steps, or unjustified leaps. The final answer should naturally and justifiably emerge from the reasoning process.<br>- Reward: The response presents a well-structured, internally consistent chain of reasoning that leads clearly and justifiably to the final answer.<br>- Penalize: The response contains contradictions, missing steps, or disconnects between reasoning and answer. Short-cut behaviors—such as giving the final answer first with unsupported or inconsistent reasoning—should also be penalized. |
| **Factuality / No Hallucination** | All claims and reasoning steps must be factually accurate and supported by the image or the prompt. The response should avoid hallucinated visual content, misidentifications, or factual inaccuracies in the reasoning process.<br>- Reward: The response is free from factual errors or hallucinations and relies only on valid observations and logical inferences.<br>- Penalize: The response introduces hallucinated details, misidentifications, or incorrect factual claims that compromise the reasoning. |
| **Reflection and Exploration** | The reasoning should demonstrate thoughtful reflection and a willingness to explore multiple possibilities, particularly when the task is ambiguous or complex. This includes acknowledging uncertainty, considering alternative interpretations, or revising initial assumptions.<br>- Reward: The response demonstrates depth through reflection, critical evaluation, or exploration of different solutions before reaching a conclusion.<br>- Penalize: The response is overly rigid, superficial, or rushed, showing little to no depth of thought, reflection, or exploration of alternative possibilities. |
| **Conciseness and Efficiency** | The reasoning should be clear, focused, and efficiently communicate the steps. It should avoid redundancy, digressions, or unnecessary elaboration that dilute the argument. For straightforward tasks, over-explaining or over-analyzing should also be avoided.<br>- Reward: The response is concise and well-structured, conveying reasoning steps precisely and proportionally to the task complexity.<br>- Penalize: The response is verbose, repetitive, or includes irrelevant content that distracts from the reasoning. It may also overthink or over-explain simple prompts. |

Table 8. Five evaluation criteria for judging *verifiable reasoning* tasks in Multi-Crit benchmark.

You are an expert in evaluating the quality of AI-generated responses according to specific evaluation criteria. Your task is to assess two responses generated by different AI assistants in reply to a user's question about an image. The image is provided as part of the input.
You must evaluate the responses **strictly and exclusively** based on the following evaluation criterion:
{Criterion}
Do not consider any other dimensions or criteria beyond what is specified above.
Here are the inputs for your evaluation:
[Question]: {Question}
[Response 1]: {Response1}
[Response 2]: {Response2}
First, provide a detailed justification for your evaluation. Refer to specific elements in the responses, how they align with the evaluation criterion, and relevant visual details from the image.
On the final line, provide your final judgment on which response is better. Your judgment must be based solely on the specified criterion. Strictly follow this format: Response X is better.

Table 9. Evaluation prompt template used for LMM judge inference in Multi-Crit. Judge models are explicitly instructed to focus strictly and exclusively on the target evaluation criterion during each inference.



Figure 8. Correlation of criterion-level human preferences for prompts that exhibit preference conflicts in the open-ended (*left*) and reasoning (*right*) splits of Multi-Crit.



Figure 9. Counts of criterion pairs exhibiting human preference conflicts in the open-ended (*left*) and reasoning splits (*right*) of Multi-Crit.

## D. Qualitative Examples

### D.1. Multimodal Prompts

Figure 10 visualizes selected multimodal prompts from Multi-Crit, spanning diverse *open-ended generation* scenarios—such as captioning, storytelling, knowledge-based tasks, and text-rich understanding—as well as *verifiable reasoning* domains including math, science, coding, and visual analogy. These prompts cover a wide range of LMM-as-a-Judge use cases.

### D.2. Full Evaluation Instances

We present six full evaluation instances—each consisting of a prompt and paired model responses—along with their multi-criterion human preference annotations. These examples illustrate how our curated data reflects fine-grained response differences and demonstrates criterion-level preference conflict patterns within each response pair.

The first three examples illustrate judgment in open-ended tasks:

- Table 10 shows an evaluation instance for judging two models in blind storytelling. This example reveals preference conflicts such as *completeness vs. grounding* and *grounding vs. no-hallucination*.

- Table 11 presents an evaluation instance for judging two models on captioning creative images. It exhibits representative preference conflicts including *completeness vs. no-hallucination*, *completeness vs. expressiveness*, and *expressiveness vs. clarity*.

- Table 12 provides an example of judging two responses from the same model on an instruction-rich image analysis task. Observed conflicts include *grounding vs. clarity* and the less common *no-hallucination vs. clarity*.

The following three examples illustrate judgment in verifiable reasoning tasks:

- Table 13 shows an evaluation instance for judging two models' thinking traces on a math reasoning task. Although both traces lead to incorrect answers, one exhibits better *visual grounding* and *explores* alternative answers, whereas the other maintains stronger *logical consistency* with *fewer perceptual hallucinations*.

- Table 14 presents an evaluation instance for judging two thinking traces produced by the same model on object counting. While both responses lead to the correct final answer, one conducts *concise* reasoning, whereas the other identifies more *visual details* and actively *reflects on each object*, but introduces *minor hallucinations* that do not affect the final result.

- Table 15 provides an example of judging two models' responses on a biological reasoning task involving a food web. In this case, one response is more *concise* and demonstrates deeper *reflection* by explicitly identifying a

key intermediate inference, while the other offers a more visually detailed description of the food-web structure.

**Fine-grained Captioning**

*Q: Describe the image in detail.*

**Semiotics & Branding**

*Q: Based on this logo alone, what can you determine about the business it represents?*

**Action Recognition/ Comic**

*Q: This is a screenshot of my employee during his work time. What is he doing?*

**Blind Storytelling**

*Q: You are a storyteller for the blind. Please tell this comic book panel by naming character and all dialogues according to the image and the following script extract.*

**Geospatial World Knowledge**

*Q: What is the best way to commute from Trump Tower to the location shown in this image?*

**OCR & Document**

*Q: Extract the text and tabular data in the image. Create a json structure with bill contents and details like date name on the bill etc.,*

**Math**

*Q: A rectangle is inscribed between $y=36-x^2$, as shown in the figure below. Write the area A of the rectangle as a function of x.*
*A: $72x-2x^3$*

**Biology**

*Q: The image shows a food web, where an arrow pointing from object A to object B indicates that A provides B with the energy needed for survival. According to the food web, what will happen if all the algae died? Choices: (A) Crabs and limpets will decrease (B) Dolphins will increase (C) Sea gulls will become extinct (D) Star fish will increase. A: A*

**Chemistry**

*Q: In the transition-state structure shown in the image, calculate the total number of bonds in the structure, including single, double, and triple bonds but excluding those involving hydrogen. A: 4*
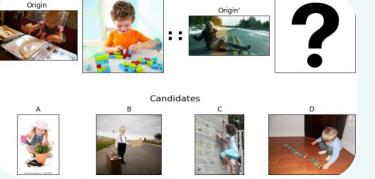
**Visual Puzzles**

*Q: Which two puzzle pieces form the larger square? Options: (A) 1 & 2 (B) 1 & 3 (C) 1 & 4 (D) 2 & 3 (E) 2 & 4 A: B*

**Visual Analogy**

*Q: Based on the relationship between Origin and Analogy, select the option from the candidates that best completes this analogy. A: B*

**Coding**

*Q: Which code snippet below can possibly create the chart in the image? We are using Python version 3.11.0, matplotlib version 3.6.3, and seaborn version 0.12.2. (A) import matplotlib.pyplot as plt... (B) ...(C)...(D)... A: B*

Figure 10. **Example prompts in the open-ended (*top*) and verifiable reasoning (*bottom*) splits of Multi-Crit.** Our benchmark spans diverse multimodal prompts, reflecting the broad evaluation domains in which LMM judges are used.

**Open-ended Example 1**



| Question | You are a storyteller for the blind. Please tell this comic book panel by naming character according to the image and the following script extract. Corresponding character's names have been added in white rectangles over the original image (do not mention rectangle or label in the storytelling). |
|---|---|
| | DIALOGUE |
| | Marco: AN INVENTOR? I DON'T BELIEVE YOU! |
| | Babi: WHY NOT? |
| Response A (llava-v1.6-vicuna-7b) | In the comic book panel, a character is named Marco. He is speaking to another character named Babi. Marco is questioning Babi's claim that he is an inventor. Babi responds with an open-ended question, encouraging Marco to reconsider his skepticism. Marco's skepticism towards Babi's claim is evident from his words, "I DON'T BELIEVE YOU!" It seems like Marco is questioning Babi's credibility or the authenticity of his claim. Babi's open-ended question, "WHY NOT?", indicates that he is not backing down and is pushing for a dialogue or a conversation where Marco can clarify his doubts. This comic panel seems to be a humorous take on skepticism and the importance of asking questions to explore different perspectives and possibilities. |
| Response B (yi-vl-plus) | In a vividly colored comic book panel, we see two characters engaging in a lively conversation. One character, with a bright yellow head and a red shirt, exclaims, "An inventor? I don't believe you!" His name, Marco, is displayed in a white rectangle. Opposite him, another character, sporting a black hairstyle and a blue shirt, confidently asks, "Why not?" His name, Babi, is also shown in a white rectangle. The contrasting colors of their clothes and hair add to the dynamic energy of their exchange, as they stand against a plain blue backdrop. The speech bubbles capture their dialogue, highlighting the playful and inquisitive spirit of their interaction. |
| Ground Truth | |

| Completeness and Coverage | A > B |
|---|---|
| Visual Grounding and Details | B > A |
| Factuality / No Hallucination | A > B |

Table 10. Example of an open-ended evaluation instance for blind storytelling. Response A conveys the moral of the story and fully follows the instruction to avoid mentioning rectangles. Response B includes more visual details but hallucinates the shirt colors.

**Open-ended Example 2**



| Question | What does this image show? |
|---|---|
| Response A (Llama-3.2-11B-Vision-Instruct) | The image depicts a surreal scene where a bottle of beer is pouring its contents onto the ground, creating a puddle that morphs into a miniature pool. Within this pool, two women are shown sitting and posing in swimsuits, as if they were enjoying a day at the beach or poolside. The juxtaposition of the beer's spillage with the serene atmosphere of the women creates a humorous and unexpected visual contrast. |
| Response B (Qwen2.5-VL-7B-Instruct) | This image is a creative and surreal depiction that combines elements of every-day objects with imaginative scenarios. It shows a large, tilted bottle pouring liquid onto the ground. The liquid spreads out in a pool around two women who are sitting on the edge of the puddle. One woman is partially submerged in the liquid, while the other sits on the edge looking at her. The scene is contrasted by the dark, plain background, which makes the bottle and the women stand out prominently. This type of imagery is often used in advertising or artistic representations to convey a message or evoke emotions. |

| Ground Truth | | |
|---|---|---|
| | Completeness and Coverage | B > A |
| | Factuality / No Hallucination | A > B |
| | Creativity and Expressiveness | A > B |
| | Clarity and Coherence | B > A |

Table 11. Example of an open-ended evaluation instance for creative-image captioning. Response A creatively highlights the humorous and unexpected effect and introduces less hallucinated content. Response B includes a brief TL;DR and an inference about image usage, with a more coherent flow across elements.
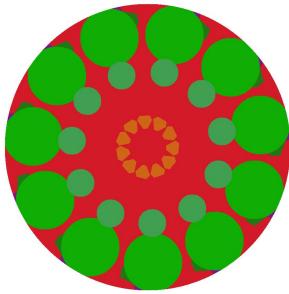
**Open-ended Example 3**



| | |
|---|---|
| Question | Create detailed one sentence image caption, capture the essence and include the following structured elements without assumptions or ambiguous language: Background and Foreground: Separate background objects and foreground objects to enhance focus; Location: The setting can be either indoors or outdoors; Shot Size: The image sizes: Extreme close up, Close up, Medium close up, Medium shot, Medium full shot, Full length. Use this parameter if a person is included in the photo; Angle: low angle shot, high angle shot, top view shot, flat lay. If the angle is straight, omit this specification; Temperature: Color temperature should be considered, with options including cold, warm, neutral, or vibrant. If the temperature is neutral, you may choose to not specify this detail; Lens: Use a fisheye lens or wide-angle as appropriate. If neither is used, this detail can be omitted; Lighting: Very important options include High-key Lighting, Low-key Lighting, Dramatic light, Backlight, silhouette, cinematic light; Composition: Employ techniques such as the rule of thirds, double exposure, over-the-shoulder shots, or symmetrical composition to structure the image; Day Time: Specify the time of day (e.g., blue hour, sunrise, morning, noon, afternoon) if this parameter is explicitly relevant; Person Caption: If a person is included, describe their race of skin color (Asian, Caucasian, Black, Arab, Indian, etc), hair (black, brown, blond, white/gray, and rarely red) view direction, emotions and mood, and clothing details. |
| Response A (GPT-4.1) | Foreground: Two children, one with light skin and straight brown hair, wearing a charcoal grey shirt and black shorts, and the other with dark skin and short curly black hair, dressed in a beige shirt and black shorts, are seated close to each other with the child on the right resting an arm around the other's shoulders; both appear calm and friendly, looking directly at the camera. Background: Metal fencing with vertical green posts is visible, along with an overcast sky and faint outlines of a streetlight or utility pole. Location: Outdoors. Shot Size: Medium shot. Angle: Low angle shot. Temperature: Neutral. Lighting: Natural diffused daylight. Composition: Centered composition, subjects positioned side by side with upper bodies and faces as focal point. Day Time: Daytime. |
| Response B (GPT-4.1) | Foreground: Two boys are sitting close together on what appears to be outdoor seating; Background: A wire mesh fence and holding poles are visible against a blank, overcast sky; Location: Outdoors; Shot Size: Medium shot; Angle: Low angle shot; Temperature: Neutral; Lighting: Diffused natural light due to an overcast day; Composition: Symmetrical composition with both boys centered in the frame; Day Time: Daytime, likely noon or afternoon given the soft daylight; Person Caption: On the left, a Caucasian boy with straight, shoulder-length brown hair faces the camera with a neutral, slightly guarded expression, wearing a charcoal grey henley shirt and black shorts; on the right, a Black boy with short, curly black hair also faces the camera with a calm, faintly smiling demeanor, dressed in a loose beige t-shirt and black shorts, his left arm draped supportively behind the other boy. |

| Ground Truth | | |
|---|---|---|
| | Completeness and Coverage | B > A |
| | Visual Grounding and Details | A > B |
| | Factuality / No Hallucination | A > B |
| | Clarity and Coherence | B > A |

Table 12. Example of an evaluation instance in instruction-rich image analysis. Both responses are randomly sampled from GPT-4.1. Response A provides more background details and is less hallucinated, while Response B better follows all required elements in the user prompt and offers a more clear and readable description of the foreground and composition.

## Verifiable Reasoning Example 1



| Question | The entire pie seen in the figure is divided among several children. Each child receives a piece of pie, and each piece of pie looks identical. The maximum possible number of children there is: Choices: (A) 10 (B) 12 (C) 9 (D) 11. |
| --- | --- |
| Answer | D |
| Response A (Qwen2.5-VL -72B-Instruct) | To determine the maximum possible number of children who can receive an identical piece of the pie, we need to analyze the symmetry of the figure. 1) Identify the Symmetry: The figure appears to be a circular design with a central red area and green circles arranged around it. The key is to find the rotational symmetry of the design. 2) Count the green circles: There are 12 green circles arranged around the central red area. This suggests that the design might have 12-fold rotational symmetry. 3) If we rotate the figure by $\frac{360°}{12} = 30°$, the figure should look the same. This means that the design is symmetric under a 12-fold rotation. 4) Conclusion: Since the design has 12-fold rotational symmetry, the pie can be divided into 12 identical pieces. Therefore, the maximum possible number of children who can each receive an identical piece of the pie is 12. Thus, the correct answer is $\boxed{B}$. |
| Response B (Gemini-2.5 -Flash) | To find the maximum possible number of children, we need to determine the number of identical pieces the pie can be divided into. This means identifying the smallest repeating unit in the circular pattern. Let's observe the elements arranged in a circle: 1) Large green circles: Count the number of large green circles arranged in the outer ring. There are 12 of them. 2) Small green/grey circles: Count the number of smaller green/grey circles in the inner ring. There are also 12 of these. 3) Orange shapes: Count the number of orange shapes in the very center. There are 12 of these as well. Since all three types of repeating elements (large green circles, small green/grey circles, and orange shapes) appear 12 times and are arranged symmetrically around the center, the pie exhibits 12-fold rotational symmetry. If each child receives an identical piece of pie, then the pie must be divided into sectors that are identical in shape and contain identical patterns. The smallest such sector would be one that contains one of each of the repeating elements. Imagine drawing radial lines from the center of the pie that equally divide the circle. If we divide the pie into 12 equal sectors (each 360/12 = 30 degrees), each sector will contain: One large green circle. One small green/grey circle. One orange shape. Since these 12 pieces would be identical and cover the entire pie, the maximum possible number of children is 12. If there were fewer children (e.g., 6), each piece would be larger and contain multiple sets of these elements (e.g., two large green circles, two small green circles, two orange shapes), but the question asks for the *maximum* number, which corresponds to the smallest possible identical pieces. Therefore, the maximum number of children is 12. The final answer is $\boxed{12}$. |
| Ground Truth | |

| | |
| --- | --- |
| Visual Grounding | B > A |
| Logical Coherence and Consistency | A > B |
| Factual Correctness / No Hallucination | A > B |
| Reflection and Exploration | B > A |
| Conciseness and Efficiency | A > B |

Table 13. Example of an evaluation instance for math reasoning. Both responses appeal to rotational symmetry but arrive at incorrect final answers. Response A is more concise and avoid hallucinated colors. Response B demonstrates slightly better grounding for the inner orange shapes and actively explores alternative solutions, yet contains some logical shifts that reduce coherence.
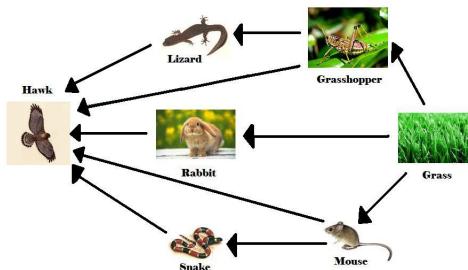
| Verifiable Reasoning Example 2 |
|---|



| | |
|---|---|
| Question | Is the number of green bus greater than the number of blue school bus? Options: (A) Yes (B) No. |
| Answer | A |
| Response A (Gemini-2.5-Flash) | 1.Identify green buses: In the image, there is one long green vehicle that is a bus. So, the number of green buses is 1. 2.Identify blue school buses: There are no blue vehicles visible in the image, and certainly no blue school buses. So, the number of blue school buses is 0. 3. Compare the numbers: 1 (green buses) is greater than 0 (blue school buses). Therefore, the number of green buses is greater than the number of blue school buses. The final answer is $\boxed{A}$ |
| Response B (Gemini-2.5-Flash) | Here's a breakdown of the vehicles in the image: - There is one green bus in the lower right portion of the image. - There is a gray dirt bike in the middle-left. - There is a golden motorcycle (possibly a chopper or custom bike) with red and green accents in the upper-middle. - There is a gray airplane with blue and green accents in the upper-right. Based on this observation: - Number of green buses = 1 - Number of blue school buses = 0 (There are no blue buses, nor any vehicles that distinctly look like a school bus of any color, let alone blue.) Since 1 is greater than 0, the number of green buses is greater than the number of blue school buses. The final answer is $\boxed{A}$ |
| Ground Truth | |

| Visual Grounding | B > A |
|---|---|
| Factual Correctness / No Hallucination | A > B |
| Reflection and Exploration | B > A |
| Conciseness and Efficiency | A > B |

Table 14. Example of an evaluation instance for object counting. Both Gemini-2.5-Flash responses correctly identify more green buses than blue ones. Response A is more compact. Response B provides finer object details and spatial relationships and explicitly checks each visible object, but introduces a minor hallucination by stating that no school bus is present, even though one appears in the image.

**Verifiable Reasoning Example 3**



| | |
|---|---|
| Question | In the diagram of the food web shown, a pointing arrow from A to B indicates that A is the source of energy that helps B survive. If the number of lizard and grasshopper decrease, the population of grass will likely? (A) decreases (B) increases (C) remains the same (D) cannot be determined from the image. |
| Answer | B |

**Response A (MiMo-VL-7B-RL)**

\<think\> I need to analyze this food web diagram to understand the relationships between organisms and determine what would happen to the grass population if lizards and grasshoppers decrease. Looking at the diagram: - Grass is eaten by: grasshopper, rabbit, and mouse - Grasshopper is eaten by: lizard and hawk - Lizard is eaten by: hawk - Rabbit is eaten by: hawk - Mouse is eaten by: hawk and snake - Snake is eaten by: hawk If the number of lizards and grasshoppers decreases: 1. Grasshoppers eat grass, so fewer grasshoppers would mean less consumption of grass. 2. Lizards don't directly eat grass, but they eat grasshoppers. With fewer lizards, there would be fewer predators controlling the grasshopper population, which could potentially lead to more grasshoppers in the long run. However, the question states that both lizard and grasshopper populations are decreasing. Since both grasshoppers (which eat grass) and lizards (which control grasshoppers) are decreasing, we have two opposing effects: - Fewer grasshoppers means less grass consumed - Fewer lizards means less predation on grasshoppers, which could lead to more grasshoppers in the future. However, since the question specifically asks about the immediate effect of these decreases, the most direct impact would be that with fewer grasshoppers consuming grass, the grass population would likely increase. \</think\> In this ecosystem: - Grass is consumed by three primary herbivores: grasshoppers, rabbits, and mice - Grasshoppers are preyed upon by both lizards and hawks - Lizards are preyed upon by hawks. When the number of lizards and grasshoppers decreases: 1. Grasshoppers directly consume grass. With fewer grasshoppers, there would be less consumption of grass. 2. Lizards don't directly consume grass, but they do control the grasshopper population through predation. With fewer lizards, there would be less predation pressure on grasshoppers, which could potentially lead to more grasshoppers in the future. However, the question specifies that both lizard and grasshopper populations are decreasing. The most immediate effect would be reduced grass consumption due to fewer grasshoppers. Since grass is also consumed by rabbits and mice (whose populations aren't mentioned as changing), the reduction in grasshoppers would likely result in increased grass availability. Therefore, the population of grass will likely $\boxed{(B) increases}$.

**Response B (Gemini-2.5-Flash)**

To determine the likely change in the grass population, we need to look at what organisms consume grass in this food web. 1. Identify consumers of Grass: The diagram shows arrows from Grass to Grasshopper, Rabbit, and Mouse. This means Grasshoppers, Rabbits, and Mice all eat Grass. 2. Analyze the given scenario: "the number of lizard and grasshopper decrease". 3. Focus on the direct impact on Grass: – The Grasshopper population decreases. Since grasshoppers eat grass, a decrease in their numbers means there will be fewer organisms consuming grass. – The Lizard population decreases. Lizards eat grasshoppers. While a decrease in lizards might lead to an increase in grasshoppers (their prey) in some scenarios, the problem explicitly states that grasshoppers *also* decrease. Therefore, the decrease in lizards doesn't counteract the effect of fewer grasshoppers eating grass. In fact, fewer lizards might be a consequence of fewer grasshoppers. 4. Conclusion: With fewer grasshoppers eating grass, the predation pressure on the grass population is reduced. Assuming other factors (like Rabbit and Mouse populations, or environmental conditions) remain constant, the grass will be consumed less. This reduced consumption will likely lead to an increase in the grass population. Therefore, if the number of lizard and grasshopper decrease, the population of grass will likely increase. The final answer is $\boxed{B}$.

| Ground Truth | | |
|---|---|---|
| Visual Grounding | | A > B |
| Reflection and Exploration | | B > A |
| Conciseness and Efficiency | | B > A |

Table 15. Example of an evaluation instance in judging biological reasoning. Both responses reach the correct conclusion (grass increases). Response A provides a more complete and detailed description of the food web. Response B is more efficient and together shows deeper reflection by explicitly recognizing that the decrease in lizards does not counteract the effect of fewer grasshoppers eating grass.

# E. Additional Experimental Results

## E.1. To Think or Not?

To study how thinking influences judge model behaviors in following diverse criteria, we examine two strong open-source LMM families that are known for producing high-quality reasoning traces across diverse domains. InternVL3.5 [39] supports both thinking (default) and non-thinking modes controlled via system prompts, enabling a clean mode-switch comparison within the same model weights. Qwen3-VL [47] provides two model variants—Instruct and Thinking—the latter specifically developed to enhance reasoning capabilities.

Table 16 presents the results, showing two trends:

1. *Smaller models benefit more from thinking, while larger models show limited gains.*
   For the smaller 8B models, both the mode-switch (InternVL3.5) and model-variant (Qwen3-VL) show clear and consistent improvements from enabling thinking, with gains across all metrics on both open-ended and reasoning splits. In contrast, the larger $\sim$30B models show more modest effects: thinking brings no significant improvement for InternVL3.5, and for Qwen3-VL it enhances reasoning judgments but slightly reduces performance on open-ended splits.
   These results suggest that explicit thinking during judgment helps smaller models better adhere to fine-grained evaluation criteria and produce criterion-specific critic reasoning. For larger models that already internalize such criteria-following capacities and reasoning patterns, additional thinking offers little benefit and may even amplify judge model biases.

2. *Thinking benefits reasoning judgments more.*
   For both Qwen3-VL and InternVL families, the thinking judge yields larger gains in evaluating verifiable reasoning than in open-ended generation tasks. Specifically, *Qwen3-VL-Thinking* outperforms its *Instruct* variant at both 8B and 32B scales, demonstrating improved ability to recognize criterion trade-offs and capture preference conflicts. This finding differs from the observation in Sec. 4.2, where RL-finetuned reasoning models on domain-specific data based on Qwen2.5-VL showed no improvement in reasoning judgments and even weakened trade-off recognition.
   This indicates that the capacity for reasoning judgment arises from general thinking abilities developed through broad reasoning training across diverse domains, rather than from narrow domain-specific fine-tuning, which often leads to overfitting and reduced generalization.

## E.2. Joint Multi-Criterion Judgment

**Task Formulation.** In our standard setting, each evaluation instance (prompt + response pair) with $K$ applicable criteria requires the LMM judge to perform $K$ separate inferences, each instantiated with a single-criterion evaluation prompt as shown in Table 9. Here, we investigate an alternative setting in which the judge performs *joint multi-criterion judgment*. For each instance, all $K$ criteria are presented to the model simultaneously, and a *single* inference is used to produce criterion-level judgments for all criteria. Table 17 shows the corresponding joint evaluation prompt: the model is instructed to treat each criterion separately while outputting all criterion-level judgments in one pass. All inference hyper-parameters remain identical to the standard single-criterion setup, allowing us to isolate the effect of joint prompting on judge model behavior.

**Results and Analysis.** We evaluate four top-performing proprietary LMMs—GPT-4o, GPT-5, o4-mini, and Gemini-2.5-Pro—under the joint multi-criterion judgment setting. The results in Table 18 reveal two major observations:

1. *Joint multi-criterion judgment affects models unevenly.*
   GPT-4o exhibits a clear performance drop across all metrics, suggesting that combining multiple criteria in a single pass amplifies its internal biases and weakens its ability to follow diverse criteria. GPT-5, in contrast, benefits from the joint setting—showing improved pluralistic accuracy and better alignment with human preferences on conflict cases. For o4-mini and Gemini-2.5-Pro, the effects of joint criterion judgment are mixed on open-ended tasks but consistently drop on reasoning domains.

2. *Joint judgment generally reduces the model's sensitivity to criterion-level trade-offs.* Across all models—except GPT-5 in the reasoning split—the trade-off sensitivity decreases. This trend is expected: given the autoregressive nature of LLMs/LMMs, generating multiple criterion-level judgments within a single inference pass inevitably induces inter-criterion dependencies, making the judge model more likely to assign the same preference direction across criteria for the same pair of responses. Table 19 shows a case where *o4-mini* fails to capture criterion-level conflicts under joint criterion judgment.

| Model | Open-ended | | | | Verifiable Reasoning | | | |
|---|---|---|---|---|---|---|---|---|
| | Pluralistic | Conflict | Tradeoff | Crit.-Avg. | Pluralistic | Conflict | Tradeoff | Crit.-Avg. |
| InternVL3.5-8B (no-think) | 23.41 | 24.95 | 44.66 | 59.09 | 27.78 | 28.83 | 54.13 | 61.96 |
| InternVL3.5-8B (think) | 25.08 | 32.34 | 62.14 | 61.04 | 32.54 | 39.15 | 69.72 | 65.85 |
| Δ | +1.67 | +7.39 | +17.48 | +1.95 | +4.76 | +10.32 | +15.59 | +3.89 |
| Qwen3-VL-8B-Instruct | 18.39 | 18.36 | 34.47 | 54.16 | 16.67 | 19.22 | 35.78 | 56.92 |
| Qwen3-VL-8B-Thinking | 24.75 | 37.33 | 66.50 | 60.61 | 38.89 | 51.25 | 83.49 | 70.70 |
| Δ | +6.36 | +18.97 | +32.03 | +6.45 | +22.22 | +32.03 | +47.71 | +13.78 |
| InternVL3.5-38B (no-think) | 29.43 | 35.93 | 61.65 | 65.36 | 40.48 | 47.69 | 75.23 | 71.37 |
| InternVL3.5-38B (think) | 30.43 | 33.73 | 64.08 | 65.10 | 37.30 | 47.69 | 75.23 | 69.82 |
| Δ | +1.00 | -2.20 | +2.43 | -0.26 | -3.18 | 0.00 | 0.00 | -1.55 |
| Qwen3-VL-32B-Instruct | 30.43 | 40.32 | 68.93 | 65.49 | 39.68 | 48.75 | 70.64 | 71.42 |
| Qwen3-VL-32B-Thinking | 29.10 | 40.12 | 67.96 | 64.88 | 43.65 | 53.38 | 80.73 | 73.88 |
| Δ | -1.33 | -0.20 | -0.97 | -0.61 | +3.97 | +4.63 | +10.09 | +2.46 |

Table 16. Comparison of thinking vs. non-thinking LMM judges. Relative improvements are shown in green and decreases in red.

You are an expert in evaluating the quality of AI-generated responses according to multiple evaluation criteria. Your task is to assess two responses generated by different AI assistants in reply to a user's question about an image. The image is provided as part of the input.

You must evaluate the responses based on the following {K} evaluation criteria. **Analyze each criterion independently and exclusively**—your judgment on one criterion should not influence your judgment on another. Do not consider any other dimensions or criteria beyond what is specified below.
{Criteria}

Here are the inputs for your evaluation:
   [Question]: {Question}
   [Response 1]: {Response1}
   [Response 2]: {Response2}

**Instructions:** For *each* criterion listed above, you must: (1) provide a detailed justification for your evaluation, referring to specific elements in the responses, how they align with that criterion, and relevant visual details from the image; and (2) on the final line of each criterion's evaluation, provide your judgment based solely on that criterion. You must choose one response as better; do not indicate a tie.
Format your response as follows for each criterion:

   **Criterion: [Criterion Name 1]**
   [Your detailed justification here]
   Judgment: Response X is better.
   **Criterion: [Criterion Name 2]**
   [Your detailed justification here]
   Judgment: Response X is better.
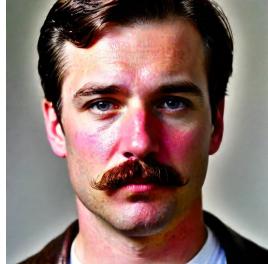   . . . (continue for all criteria)

Strictly follow this format: ``Response X is better.'' on the last line of each criterion block.

Table 17. Joint multi-criterion evaluation prompt for LMM judges. All applicable criteria for an evaluation instance are assessed in a single pass, and the judge is explicitly instructed to treat each criterion separately.

| Model | | Open-ended | | | Verifiable Reasoning | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Pluralistic | Conflict | Tradeoff | Criterion Avg. | Pluralistic | Conflict | Tradeoff | Criterion Avg. |
| GPT-4o | 31.44 | 44.91 | 66.02 | 69.57 | 41.27 | 55.16 | 84.40 | 69.79 |
| + Joint multi-criterion | 30.10 | 26.35 | 38.83 | 64.91 | 34.13 | 34.52 | 58.72 | 63.15 |
| Δ | -1.34 | -18.56 | -27.19 | -4.66 | -7.14 | -20.64 | -25.68 | -6.64 |
| GPT-5 | 29.77 | 38.52 | 62.62 | 68.51 | 45.24 | 56.58 | 78.90 | 77.41 |
| + Joint multi-criterion | 34.78 | 43.51 | 60.68 | 70.43 | 47.62 | 61.57 | 82.57 | 79.58 |
| Δ | +5.01 | +4.99 | -1.94 | +1.92 | +2.38 | +4.99 | +3.67 | +2.17 |
| o4-mini | 32.78 | 43.11 | 64.56 | 69.67 | 53.17 | 65.84 | 83.49 | 80.85 |
| + Joint multi-criterion | 36.79 | 40.12 | 58.74 | 68.82 | 45.24 | 55.87 | 77.06 | 75.61 |
| Δ | +4.01 | -2.99 | -5.82 | -0.85 | -7.93 | -9.97 | -6.43 | -5.24 |
| Gemini-2.5-Pro | 28.76 | 37.92 | 66.50 | 63.67 | 41.27 | 52.33 | 75.93 | 73.06 |
| + Joint multi-criterion | 32.11 | 41.12 | 62.62 | 67.06 | 35.71 | 40.21 | 56.88 | 70.96 |
| Δ | +3.35 | +3.20 | -3.88 | +3.39 | -5.56 | -12.12 | -19.05 | -2.10 |

Table 18. Comparison of joint multi-criterion judgment against standard single-criterion judgment. Relative improvements are shown in green, decreases in red.

| | |
|---|---|
| Case Study: Joint Multi-Criterion Failure of the *o4-mini* Judge in Capturing Preference Conflicts | |



| | |
|---|---|
| Question | Provide a single 10-year age bracket (e.g., "30s," "40s," etc.) and describe the man's hair color and hairstyle based only on the image. |
| Response 1 | Based on the image provided, the man appears to be in his 30s. His hair is dark brown and is styled in a classic fashion, swept to the side, which complements his period attire. The style is reminiscent of early 20th-century fashion, which aligns with the vintage feel of the image. |
| Response 2 | The man in the photo appears to be in his late 30s to early 40s. His hair is dark brown with a slight wave pattern. It is neatly styled, suggesting that he takes care of his appearance. The hair is cut short, and the parting on the side complements his overall look. |
| Ground Truth | Completeness and Coverage     Resp. 1 > Resp. 2<br>Factuality / No Hallucination     Resp. 2 > Resp. 1 |
| Single-Criterion Judgment (*o4-mini*) | **Completeness and Coverage**<br>The prompt asks for a single 10-year age bracket (e.g., "30s," "40s," etc.). Response 1 clearly states "in his 30s," satisfying that requirement exactly. Response 2 instead gives "late 30s to early 40s," which spans two decades and thus does not fulfill the prompt's decade-long range specification. - Both responses describe the hair color ("dark brown") and hairstyle. However, Response 1 stays focused on the core question by succinctly noting the side-swept, classic style evident in the image. Response 2 adds extra inferences ("suggesting that he takes care of his appearance") and splits the age range improperly, making it less precise with respect to the user's request. Because Response 1 fully and directly addresses both parts of the question—providing the correct decade-long age bracket plus hair color and style—while Response 2 does not, Response 1 is the better answer under the completeness and coverage criterion.<br>Response 1 is better. ✓<br><br>**Factuality / No Hallucination**<br>Response 1 introduces unsupported details about "period attire" and "early 20th-century fashion," neither of which are evident in the image—this is a visual hallucination. Response 2, by contrast, sticks to what is visible (age estimate spanning late 30s to early 40s, dark brown hair, slight wave, side part, neatly styled) without inventing extraneous context.<br>Response 2 is better. ✓ |
| Joint Multi-Criterion Judgment (*o4-mini*) | **Criterion: Completeness and Coverage** Response 1 directly answers both parts of the prompt within the requested format: it gives a single 10-year range ("30s") and describes hair color ("dark brown") and style ("classic fashion, swept to the side"), even tying in contextual cues ("period attire," "vintage feel"). Response 2, by contrast, straddles two decades ("late 30s to early 40s"), which fails to provide a single 10-year band as requested, and although it describes hair color and neat styling, it does not anchor its observations as neatly to the prompt's format. Judgment: Response 1 is better. ✓<br>**Criterion: Factuality / No Hallucination** Both responses correctly observe dark brown hair and a side part, but Response 2 introduces imprecision by spanning "late 30s to early 40s," which contradicts the user's instruction for a single decade. Response 1 sticks to a plausible age ("30s") based on the man's appearance. Neither response hallucinates details about color, but Response 2's age estimate violates the prompt's constraints. Judgment: Response 1 is better. ✗ |

Table 19. **In this case, *o4-mini* correctly identifies the *completeness vs. no-hallucination* conflict under single-criterion judgments but collapses to a uniform preference under joint multi-criterion judgment.** With single-criterion evaluation prompts, the model successfully adheres to each criterion independently and captures the preference disagreement, whereas under joint prompting, completeness dominates and the model fails to follow the no-hallucination criterion.