

# MULTI-OBJECTIVE TASK-AWARE PREDICTOR FOR IMAGE-TEXT ALIGNMENT

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Evaluating image-text alignment while reflecting human preferences across multiple aspects is a significant issue for the development of reliable vision-language applications. It becomes especially crucial in real-world scenarios where multiple valid descriptions exist depending on contexts or user needs. However, research progress is hindered by the lack of comprehensive benchmarks and existing evaluation predictors lacking at least one of these key properties: (1) *Alignment with human judgments*, (2) *Long-sequence processing*, (3) *Inference efficiency*, and (4) *Applicability to multi-objective scoring*. To address these challenges, we propose a plug-and-play architecture to build a robust predictor, MULTI-TAP (**M**ulti-**O**bjective **T**ask-**A**ware **P**redictor), capable of both multi and single-objective scoring. MULTI-TAP can produce a single overall score, utilizing a reward head built on top of a large vision-language model (LVLMs). We show that MULTI-TAP is robust in terms of application to different LVLM architectures, achieving significantly higher performance than existing metrics (e.g., +42.3 Kendall’s  $\tau_c$  compared to IXCREW-S on FlickrExp) and even on par with the GPT-4o-based predictor, G-VEval, with a smaller size (7–8B). By training a lightweight ridge regression layer on the frozen hidden states of a pre-trained LVLM, MULTI-TAP can produce fine-grained scores for multiple human-interpretable objectives. MULTI-TAP performs better than VisionREWARD, a high-performing multi-objective reward model, in both performance and efficiency on multi-objective benchmarks and our newly released text-image-to-text dataset, EYE4ALL. Our new dataset, consisting of chosen/rejected human preferences (EYE4ALLPref) and human-annotated fine-grained scores across seven dimensions (EYE4ALLMulti), can serve as a foundation for developing more accessible AI systems by capturing the underlying preferences of users, including blind and low-vision (BLV) individuals. Our contributions can guide future research for developing human-aligned predictors.

## 1 INTRODUCTION

Accurate and efficient evaluation of image-text alignment is a fundamental task in multimodal research, serving as a key benchmark for assessing large vision-language models (LVLMs) (Lin et al., 2014; Hossain et al., 2019; Ghandi et al., 2023). As LVLMs are deployed in complex real-world scenarios, such as assistive technologies (Bandukda et al., 2019; Kazemi et al., 2023; Kuriakose et al., 2023; Chidiac et al., 2024) and instructional agents (Wang et al., 2024c; Li et al., 2024), the demand for human-aligned evaluation protocols for multimodal inputs has significantly increased (e.g., automatic metrics (Grimal et al., 2024; Hartwig et al., 2024) and alignment training (Christiano et al., 2017; Schulman et al., 2017; Ahmadian et al., 2024)). Existing model-based image-text alignment predictors, also referred to as model-based metrics or reward models, can be categorized into three types: (a) encoder-based predictors, (b) text-based scoring predictors with generative LVLMs (generative reward models), and (c) scalar-based scoring predictors with generative LVLMs (scalar-based reward models). While each shows distinct strengths, none simultaneously satisfies four key properties: (1) *Strong correlation with human judgments*, (2) *Long-sequence processing*, (3) *Inference efficiency*, and (4) *Applicability to multi-objective scoring*. Meeting all four is essential for capturing diverse and context-dependent user preferences.

For example, (a) encoder-based predictors such as CLIP-Score (CLIP-S) (Hessel et al., 2021), BLIP-Score (BLIP-S) (Li et al., 2022), and related variants (Xu et al., 2024b; Sarto et al., 2023; Wada et al.,

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

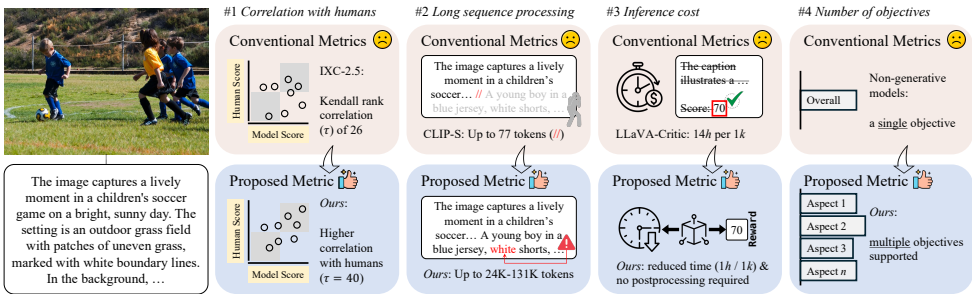


Figure 1: Comparison between existing image-text alignment metrics and ours. Our proposed method, applicable to different types of LVLMs, overcomes the challenges of conventional metrics in terms of (1) showing high correlations with human judgments, (2) understanding long input text sequences with detailed instructions, (3) reducing the inference time by returning precise scalar-based scores, and (4) enabling interpretable embeddings disentangled to multi-objective scores.

2024; An et al., 2024) are lightweight and efficient, yet their inherently limited context windows make understanding long text challenging (#2 in Figure 1). On the other hand, (b) generative reward models (Deitke et al., 2024; Xiong et al., 2024; Meta, 2024; Chen et al., 2024; Wang et al., 2024b) show improved semantic understanding supported by larger input size and alignment training. While effective, they are comparatively more computationally intensive (#3 in Figure 1) and often require prompt-tuning and additional post-processing (Li et al., 2024; Tong et al., 2024; Lambert et al., 2024), including bias subtraction to prevent ranking distortions (Zhu et al., 2025).

Finally, the most recent open-source scalar-based predictor, exemplified by InternLM-XComposer-2.5-Reward (IXCREW-S) (Zang et al., 2025), effectively addresses the aforementioned problems. Specifically, simply appending a scoring head on top of InternLM-XComposer-2.5 (InternLM) (Zhang et al., 2024c), it significantly improves cost efficiency and shows robust performance on the general vision-language reward benchmarks (Li et al., 2024). However, we observe weak agreement with the human judgments, captured by significantly low Kendall’s  $\tau$  on several image-text alignment benchmarks (Xu et al., 2019; Plummer et al., 2015) (#1 in Figure 1). Moreover, the model is released to be tied to a single LVLm backbone, InternLM, which constrains architectural modularity.

To address these limitations, we introduce a novel, LVLm-based predictor called MULTI-TAP (Multi-Objective Task-Aware Predictor), capable of producing human-aligned scores for both single and multiple-objective across diverse criteria (#4 in Figure 1). MULTI-TAP produces a robust single overall score and fine-grained scores aligned with multiple human-interpretable dimensions, utilizing the last hidden states from our newly trained reward model built on top of LVLm. Our predictor outperforms VisionReward (VisionREW-S) (Xu et al., 2024a), the only publicly available multi-objective scalar-based reward model for multimodal input (Xu et al., 2024a; Team, 2024). Here, to avoid confusion, “VisionREW-S/ImgREW-S” denotes models and “VisionREW/ImgREW” denotes datasets. Importantly, we instantiate the framework with widely used LVLms, including Qwen2-VL (Wang et al., 2024b), InternLM (Zhang et al., 2024c), and LLaMA-3.2 (Meta, 2024), demonstrating model agnosticity, and at the same time, tackling all four core challenges in human-aligned evaluation.

To further validate our approach in a more challenging and practical setting, we introduce a novel text-image-to-text (TI2T) dataset, EYE4ALL, built upon judgments of 25 human annotators, including crucial perspectives from the blind and low-vision (BLV) individuals. Unlike existing datasets that focused on evaluating the quality of generated images (Xu et al., 2024a; Zhang et al., 2024e), EYE4ALL contains human judgments on the quality of the LVLm-generated text response and alignment to text request and scenery image. This unique BLV-centered benchmark is carefully curated, inspired by the BLV preference analysis from An et al. (2025). Specifically, the human annotators are guided to evaluate responses with respect to the BLV perspectives (given the BLV-driven request) rather than merely verifying the consistency of various image responses. For multiple evaluation purposes, we provide two complementary modes: EYE4ALLPref, consisting of the human preferences on two different LVLm text responses, and EYE4ALLMulti, a collection of human judgment scores across fine-grained dimensions, such as accuracy, sufficiency, and safety, facilitating both single- and multi-objective scoring evaluations. Our EYE4ALL covers diverse

pedestrian scenarios, enabling comprehensive assessment and systematic evaluation of recent LVLMs within realistic navigation contexts.

In summary, our study makes the following contributions: (1) **Scalable single- and multi-objective reward modeling framework on LVLM** for developing a robust scalar-based human-aligned predictor on multimodal datasets. (2) **MULTI-TAP**, strongly aligned with human judgments on both single and multiple dimensions. (3) **EYE4ALL**, a response-quality-oriented benchmark designed for practical evaluation and building robust assistive AI systems. Our work can serve an important role in guiding future research on robust and human-centered multimodal evaluation.

## 2 RELATED WORKS

### 2.1 IMAGE-TEXT ALIGNMENT EVALUATION PREDICTORS

Model-based metrics aim to automatically assess image-text alignment by approximating human judgment. The most widely used predictors are encoder-based metrics, which are commonly divided into *reference-based* and *reference-free*, where a reference denotes a human-written ground-truth caption paired with an image. Reference-based metrics (e.g., Polos (Wada et al., 2024) and RefPAC-S (Sarto et al., 2024)) score the alignment between an image and a candidate caption conditioned on one or more references. In contrast, reference-free metrics (e.g., CLIP-S (Hessel et al., 2021), BLIP-S (Li et al., 2022)) score the alignment based solely on image and candidate caption pair. Reference-based methods generally correlate more strongly with human ratings, but they require costly reference annotations. Inspired by the LLM-as-a-judge concept, recent work has explored generative LVLMs for evaluation (Deitke et al., 2024; Xiong et al., 2024; Meta, 2024). For instance, Tong et al. (2024) utilizes GPT-4o (OpenAI, 2024b) with the Chain-of-Thoughts (CoT) reasoning prompt to evaluate the alignment. In parallel, scalar-based reward models have also emerged by training a linear projection head on top of generative LVLMs (Zang et al., 2025). Nevertheless, existing multi-objective reward models, such as VisionREW-S (Xu et al., 2024a) and MPS (Zhang et al., 2024e), suffer from limited efficiency and accessibility.

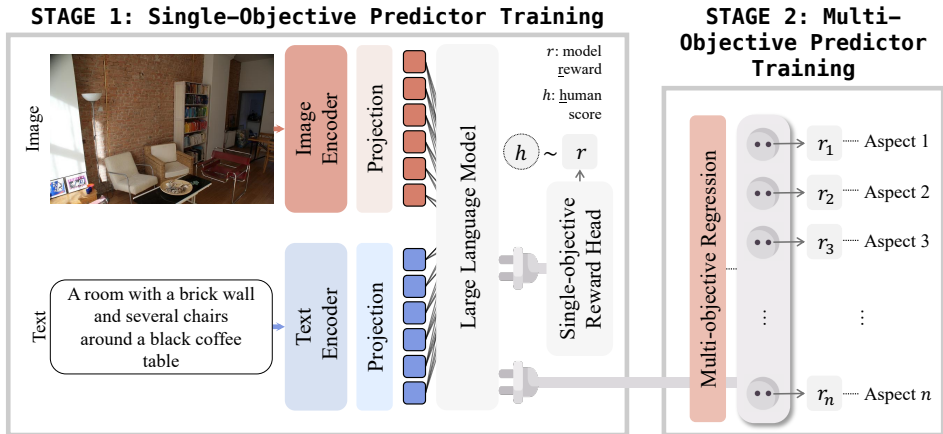
### 2.2 IMAGE-TEXT ALIGNMENT EVALUATION DATASETS

Datasets for image-text alignment evaluation are typically categorized by two annotation settings: *pointwise* and *pairwise*. For the pointwise ranking datasets, each sample is labeled with the absolute human judgment scores across fine-grained scales (Xu et al., 2024b; Wada et al., 2024; Plummer et al., 2015), enabling nuanced metric analyses. Evaluations on these datasets often report Kendall’s  $\tau$  correlation to quantify alignment between metric scores and human judgments. Pairwise ranking datasets, in contrast, provide preference labels between two competing candidates (Xu et al., 2019; Shekhar et al., 2017). Instead of having absolute scores as the labels, the datasets are labeled as positive (human-preferred) or negative (human-rejected). Several pointwise ranking datasets, such as OID (Krasin et al., 2017) and Polaris (Wada et al., 2024), include multiple candidate captions with human judgment scores or ground-truth captions, enabling them to be repurposed for either pointwise or pairwise evaluations (An et al., 2025). Despite recent advances, few multimodal datasets are labeled with human judgment scores depending on varying preferences of users and criteria (Xu et al., 2024a; Team, 2024; An et al., 2025; Kang et al., 2025). This dataset scarcity hinders the development of LVLMs that can adapt their responses to varying contexts and user needs.

## 3 MULTI-OBJECTIVE TASK-AWARE PREDICTOR (MULTI-TAP)

We present MULTI-TAP (Figure 2), a robust image-text alignment predictor that supports single- and multiple-objective scoring. Our predictor can return to output either an overall score or multiple scores across their uniquely defined criteria. Training proceeds in two stages. At Stage 1, we train a single-objective predictor to produce a single, unified score that captures overall semantic alignment between images and texts, while simultaneously shaping semantically rich multimodal embeddings for Stage 2. During Stage 2, we use these frozen embeddings to build a *multi-objective, task-aware* predictor that produces scores across multiple human-interpretable dimensions. **For inference, stage 1 is required for training to yield an overall image-text alignment score, while training both stages returns multiple scores across human-interpretable, fine-grained dimensions.** To the best of our

162 knowledge, MULTI-TAP is the first human-aligned reward modeling framework explicitly designed  
 163 for image-text alignment in accessibility-critical contexts.  
 164



165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180 **Figure 2: Schematic diagram of proposed MULTI-TAP architecture.** At Stage 1, MULTI-TAP produces a scalar value reflecting image-text alignment by appending a reward head to the LVLM. For Stage 2, a ridge regression layer is added to the trained multimodal embeddings, generating scores across multiple aspects. [We train the single-objective reward head and multi-objective regression layer in Stages 1 and 2, respectively.](#)

184  
185 **3.1 SINGLE-OBJECTIVE PREDICTOR TRAINING**

186 As shown in Figure 2, our architecture appends a reward head to the generative LVLM. This design  
 187 allows the LVLM to process extended inputs and generate a unified semantically rich multimodal  
 188 representation (#2 in Figure 1). The reward head maps the last hidden states of LVLM, or multimodal  
 189 embeddings, to a *scalar* score, adaptable for single-objective scoring. It significantly reduces the  
 190 inference latency compared to text-based generative scoring methods (#3 in Figure 1). We initialize  
 191 the reward head using a zero-centered Gaussian distribution with standard deviation  $\frac{1}{\sqrt{d+1}}$  ( $d$ : hidden  
 192 dimension), following the standard initialization practices (von Werra et al., 2020). Departing from  
 193 the Bradley-Terry (BT) model-based losses (Bradley & Terry, 1952) generally used in reward model  
 194 training (Christiano et al., 2017; Stiennon et al., 2020), we adopt mean squared error (MSE) loss to  
 195 explicitly align score outputs with human judgment scores:  $\min_{\theta} \sum_i^N (r_i - h_i)^2$ , where  $\theta$  denotes  
 196 parameters of LVLM and reward head,  $r$  is the predicted scalar score,  $N$  is the number of samples,  
 197 and  $h$  is the human judgment score. Our empirical findings indicate that MSE, with its convex  
 198 formulation and simpler optimization, offers superior training stability and performance [compared to the BT loss, where unstable spikes are frequently observed.](#) We train LVLM and a reward model on  
 199 two publicly available image-text alignment datasets: Polaris (Wada et al., 2024) and ImageReward  
 200 (ImgREW) (Xu et al., 2024b), allowing the comprehensive alignment training of text quality judgment  
 201 in terms of images and vice versa (details in Section 5). Training on these datasets helps the model  
 202 to produce a robust score aligned with human judgment (#1 in Figure 1) and to create meaningful  
 203 multimodal embeddings for multi-objective scoring (#4 in Figure 1).  
 204

205  
206 **3.2 MULTI-OBJECTIVE PREDICTOR TRAINING**

207 Optimizing predictors directly under multi-objective settings is challenging. Prior work has shown  
 208 that gradient-based optimization methods (e.g., MGDA (Zhang et al., 2024d), PCGrad (Yu et al.,  
 209 2020)) are costly at the scale of LVLMs and often fail to balance inherently conflicting objectives  
 210 (He & Maghsudi, 2025). Ensemble-based approaches such as EMORL (Kong et al., 2025) partially  
 211 alleviate these difficulties by aggregating models trained on individual objectives, but they still face  
 212 instability and limited interpretability ([without post-hoc calibration](#)) when deployed at scale. These  
 213 observations suggest that directly training a multi-objective predictor with LVLM is both expensive  
 214 and unreliable.

215 To address this, we adopt a two-stage training paradigm. We reuse the multimodal embeddings  
 from the previous stage by appending a multi-objective regression layer to the model, excluding

the original reward head. Inspired by previous work (Wang et al., 2024a), we treat the LVLM as a feature extractor (the plug symbol in Figure 2). Let  $z_i \in \mathbb{R}^d$  denote the frozen multimodal embedding and let  $y_i \in \mathbb{R}^K$  be the vector of human scores for  $K$  dimensions. A ridge regression head predicts  $\hat{y}_i = Wz_i + b$  with parameters  $W \in \mathbb{R}^{K \times d}$  and  $b \in \mathbb{R}^K$  trained by  $\min_{W,b} \sum_{i=1}^N \|y_i - Wz_i - b\|_2^2 + \alpha \|W\|_F^2$ . This head outputs multiple scalar scores per sample, each aligned with a human-interpretable criterion, improving transparency and supporting efficient customization (#4 in Figure 1). Utilizing precomputed hidden states from the frozen LVLM backbone, these regression heads can be trained asynchronously to meet task-specific needs. We opt for ridge regression as a principled design choice to maximize model interpretability, efficiency, and scalability, especially in deployment scenarios where full finetuning or gradient-based optimization is prohibitive.

Unlike ArmoRM and VisionREW-S (Xu et al., 2024a), which aggregate multi-objective outputs into a single score through a learned linear or gating head, we do not aggregate multi-objective outputs into a single overall score. This is because score aggregation across predefined dimensions may fail to holistically capture overall quality. For instance, VisionREW-S has several zero values in its aggregation head, which means the weights for specific dimensions are unstable and uninformative for computing an overall score. Instead, we use a single overall reward score from Stage 1, and a set of dimension-specific scores derived via lightweight ridge regression training from Stage 2.

## 4 EYE4ALL: BLV-AWARE, MULTI-OBJECTIVE IMAGE-TEXT ALIGNMENT SCORING DATASET

To supplement the limited pool of current multi-objective scoring datasets containing fine-grained human judgment scores for image-text alignment evaluation and to validate our MULTI-TAP in a more realistic setting, we introduce EYE4ALL, a curated evaluation dataset. This section outlines the LVLM generation pipeline and the human annotation protocol (detailed steps are in Appendix A).

### 4.1 LVLM RESPONSE COLLECTION

We first collect diverse text responses from LVLMs conditioned on an image and a scene-relevant request. Existing multi-objective image-text alignment scoring datasets (Xu et al., 2024a; Team, 2024; Zhang et al., 2024e) typically rely on judgment scores concerning image quality and emphasize generic preferences. In contrast, EYE4ALL consists of reasoning chains that are useful and applicable to (but not limited to) Blind and Low-Vision (BLV) users for navigational purposes in daily lives. We use Sideguide (Park et al., 2020) and Sidewalk (AIHub, 2019) scenery image corpora, pairing each image with BLV-plausible textual requests aligned to the depicted scene (An et al., 2025).

We collect text responses by prompting QWEN2-VL (Wang et al., 2024b), LLaVA-1.6 (Liu et al., 2023), and InternLM-XComposer2-VL (InternLM-X2-VL) (Dong et al., 2024), which are known to demonstrate remarkable in-context learning ability (Zong et al., 2025). These models are instructed to generate responses from the perspective of BLV users, utilizing the model instructions and the BLV-plausible requests from An et al. (2025). For each instance, we randomly select one of the three different LVLM responses and refine the responses by prompting GPT-4o mini (OpenAI, 2024b) the original LVLM responses (system and few-shot prompts are in Appendix B). **This process of using GPT-4o to refine responses from various models rather than generating them from scratch helps to minimize the introduction of strong bias induced by a single model.**

### 4.2 HUMAN ASSESSMENT FOR EVALUATING LVLM RESPONSES

We collect human judgments from 25 sighted human annotators (approved by the Institutional Review Board), evaluating the collected LVLM responses in terms of two main aspects: (1) whether the text response aligns with the image and the text request, and (2) whether the response addresses potential safety concerns evident in the scenery image that BLV users might otherwise miss. To balance efficiency and diversity of annotations, we randomly sample 1k image-request-response triplets from the previous stage. Each annotator completes 100 items in an online setting, typically taking 2 to 3 hours.

For each sample, the annotators are required to rate the refined GPT-4o responses along seven dimensions: (1) *Direction Accuracy*, (2) *Depth Accuracy*, (3) *Safety*, (4) *Sufficiency*, (5) *Conciseness*, (6) *Hallucination*, and (7) *Overall Quality*. These criteria were selected based on the importance and challenges of producing safe and informative LVLM responses given the image and request (Karamolegkou et al., 2025). The accuracy criterion is divided into two aspects (*Direction* and *Depth*) since precise spatial guidance directly affects user safety, particularly for BLV users. In addition, we exclude requests that contain horizontal directions beyond the 9 to 3 o’clock, since directions from 4 to 8 o’clock correspond to areas behind the viewer and are not visible in standard non-panoramic images. Unlike the other criteria, the *Hallucination* is assessed in a dichotomous format to capture misleading or false navigational content. This is distinct from *Safety*, where the focus is on whether the response includes obstacles or provides safety-relevant guidance.

All aspects are annotated with an averaged scalar score from 2–3 human judgment scores on a Likert scale of 1 to 5, except for the *Hallucination*, which is evaluated as either 0 (presence) or 1 (absence), depending on these three issues: (1) non-related information, (2) inaccurate step-by-step order, and (3) repeated content. The other six dimensions are evaluated with a fine-grained rubric (details in Appendix B). The human annotation procedure yields 2,112 unique samples that constitute the EYE4ALLMulti dataset for multi-objective scoring assessment. In addition, we also collect high-quality human-refined captions for each image (and request)-response pair that satisfy all seven evaluation criteria, labeled as positive (preferred) samples in our human preference dataset, EYE4ALLPref. The well-constructed predictors should understand the scenery images as well as context-dependent requests of our EYE4ALL to simulate human judgment preferences and fine-grained scores across multiple dimensions (more details are in Appendix B).

## 5 EXPERIMENTS

We evaluate MULTI-TAP through comprehensive experiments spanning both single- and multi-objective settings. Our study benchmarks against a wide range of predictors and datasets, enabling a rigorous assessment of its robustness and efficiency across diverse tasks and model architectures.

### 5.1 MULTI-TAP TRAINING

**Datasets.** We train the single-objective predictor using two open-sourced datasets with complementary annotation protocol: Polaris (Wada et al., 2024) and ImageReward (ImgREW) (Xu et al., 2024b). Polaris is a dataset containing annotations of captions based on given images; on the other hand, ImgREW is a dataset evaluating the quality of generated images according to prompts. The Polaris dataset consists of general image-caption alignment scores ranging from 0 to 1 (discretized into 0.0, 0.25, 0.5, 0.75, 1.0). ImgREW also consists of scores, ranging from 1–7 (later normalized to 0–1), which measure prompt-image alignment. *Note that these two datasets are all open-source, and we adhere to the predefined training/test splits, following standard practice.*

**Models.** To examine the generalization across model scales and architectures, we instantiate MULTI-TAP on the following widely used LVLMs: Qwen2-VL-2B/7B (Wang et al., 2024b), InternLM-XComposer-2.5-7B (InternLM-7B) (Zhang et al., 2024c), and LLaMA-3.2-11B (Meta, 2024). For the 2B model, we set the learning rate to  $2e-7$ , and for the larger models, we use a learning rate of  $2e-6$ . We utilize 8 A100 GPUs for training LLaMA-3.2-11B and 8 RTX A6000 GPUs for the others with seed 42. All models are trained for a single epoch, with a batch size of 8 and a gradient accumulation of 4. For the multi-objective version of MULTI-TAP, the backbone LVLM is frozen, and only a lightweight ridge regression head over the final hidden states is trained. We perform a hyperparameter search of the regularization coefficient  $\alpha$  within the scope of  $\{0.001, 0.01, 0.1, 1, 10, 100\}$ , selected based on the lowest training loss. We train MULTI-TAP with a single epoch for each  $\alpha$  *since we empirically observe rapid and stable convergence during the first epoch (for both stages).*

### 5.2 MULTI-TAP EVALUATION

**Datasets.** We evaluate the efficacy of MULTI-TAP across diverse image-text alignment benchmarks: PASCAL-50S (Xu et al., 2019), FOILR1/R4 (R1 and R4 refer to evaluation using one and

	Pairwise Ranking Datasets					Pointwise Ranking Datasets			
	PASCAL P-Acc	FOILR1 P-Acc	FOILR4 P-Acc	Polaris* P-Acc	OID* P-Acc	ImgREW P-Acc	FlickrExp $\tau_c$	FlickrCF $\tau_b$	Polaris $\tau_c$
<i>CLIP/BLIP-based predictors</i>									
CLIP-S	80.7	87.2	87.2	79.7	56.5	56.7	51.2	34.4	52.3
LongCLIP-S	82.8	91.6	91.6	77.5	58.1	56.5	54.1	35.4	54.0
PAC-S	82.4	93.7	94.9	77.0	57.7	57.2	55.9	37.6	52.5
Ref-free Polos	81.0	88.7	88.7	60.0	66.2	56.6	51.4	34.4	52.3
RefCLIP-S†	83.1	91.0	92.6	-	-	-	53.0	36.4	52.3
RefPAC-S†	84.7	88.7	94.9	-	-	-	55.9	37.6	56.0
Polos†	86.5	93.3	95.4	-	-	-	56.4	37.8	57.8
BLIP-S	82.5	95.1	95.1	79.5	59.3	57.8	57.1	37.8	54.0
ImgREW-S	81.5	93.8	93.8	73.3	58.5	65.2	49.8	36.2	52.3
<i>Reward model-based predictors</i>									
IXCREW-S	74.2	94.3	94.3	81.9	57.5	53.6	17.0	25.7	50.5
<b>MULTI-TAP</b>									
- Qwen-2B-S	81.5	<b>98.0</b>	<b>98.0</b>	<b>87.0</b>	53.2	59.0	56.8	38.9	60.1
- Qwen-7B-S	84.0	97.8	97.8	82.2	58.1	63.2	58.1	38.5	61.1
- InternLM-7B-S	83.2	96.5	96.5	81.6	61.0	61.6	<b>59.3</b>	<b>39.5</b>	<b>61.6</b>
- LLaMA-3.2-S	83.0	96.9	96.9	78.8	<b>68.7</b>	62.2	56.8	38.0	60.7

Table 1: Performances of various predictors (S: Score) on image-text alignment datasets. For both pairwise and pointwise ranking evaluation, MULTI-TAP models consistently outperform other metrics. Reference-based metrics, marked with †, cannot be evaluated on datasets without references (indicated by “-”).

four references) (Shekhar et al., 2017), Polaris\* (Wada et al., 2024), OID\* (Krasin et al., 2017), ImageReward (ImgREW) (Xu et al., 2024b), Flickr8k-Expert (FlickrExp) and Flickr8k-CF (FlickrCF) (Plummer et al., 2015), and Polaris (Wada et al., 2024). An asterisk (\*) indicates that the original dataset has been reformulated into a pairwise comparison format by binarizing scores at the median threshold to separate preferred from rejected samples. For OID\*, we use a curated 246-sample subset of exact matches due to partial availability. The first five pairwise ranking datasets are evaluated with pairwise accuracy (P-Acc), where a higher score for the positive sample indicates a correct answer. The pointwise ranking datasets are evaluated with Kendall’s correlation coefficients ( $\tau_b$  or  $\tau_c$ ), scaled by 100 for comparability with accuracy. To further evaluate performance on long-form and diverse prompts (image-to-text; I2T, text-to-image; T2I, and text-image-to-text; TI2T), we also test on Sighting (Kang et al., 2025), Align-anything (Team, 2024), and our EYE4ALLP ref.

For multi-objective scoring evaluation, we benchmark predictors on VisionREW (Xu et al., 2024a), Align-anything (Team, 2024), and our EYE4ALLMulti datasets. Since the test set of VisionREW is not publicly available, we randomly sample 1k samples from its training data for evaluation and use the remainder for training. For EYE4ALLMulti, we construct an additional 1k training set comprising scene-request pairs and responses of GPT-4o mini (OpenAI, 2024b). This dataset results in a diverse score range for each criterion, avoiding model overfitting for one particular dimension. Since VisionREW-S only outputs binary-scaled scores, we apply median-based binarization to Align-anything (TI2T-Binary, T2I-Binary) and EYE4ALLMulti (e.g., using a threshold of 2 on a 1–4 scale). MULTI-TAP, in contrast, outputs continuous scores; we report the uncalibrated multi-objective scores as the baselines (see Figure 10 in Appendix C).

**Comparison Baselines.** In the single-objective setting, we compare MULTI-TAP with a broad spectrum of predictors. First, we include the CLIP and BLIP-based metrics: CLIP-S/RefCLIP-S (Hessel et al., 2021), LongCLIP-S (Zhang et al., 2024a), PAC-S/RefPAC-S (Sarto et al., 2023), Ref-free Polos/Polos (Wada et al., 2024), BLIP-S (Li et al., 2022), and ImageReward (ImgREW-S) (Xu et al., 2024b). We also report reference-based metrics run in a reference-free configuration. In addition, we include generative reward models and scalar-based reward models: Molmo-7B (Deitke et al., 2024), LLaVA-Critic-7B (Xiong et al., 2024), Qwen2-VL-7B (Wang et al., 2024b), InternVL2-8B (Chen et al., 2024), and IXCREW-S (Zang et al., 2025). Lastly, we compare ours with G-VEval (Tong et al., 2024), which leverages GPT-4o mini (OpenAI, 2024b). To compare MULTI-TAP with existing multi-objective models, we select VisionREW-S (Zhang et al., 2024b), the only publicly available multimodal multi-objective scalar-based predictor to our knowledge. We employ the BF16 release with default settings and disable the masking method for better performance in our runs.

## 6 RESULTS

### 6.1 CORRELATION WITH HUMAN JUDGMENTS

We first show that our proposed MULTI-TAP models generally align well with human judgments in terms of *both* pairwise and pointwise ranking datasets. As shown in Table 1, four versions of MULTI-TAP, built on different LVLMs, mainly outperform conventional CLIP-, BLIP-, and reward model-based predictors across a wide range of image-text alignment benchmarks. Regardless of the architectures, MULTI-TAP generally achieves higher performances than the existing scalar-based reward model, IXCREW-S. In particular, MULTI-TAP<sub>Qwen-2B-S</sub> notably shows the best accuracy performances on FOIL and Polaris\*, achieving 98.0% and 87.0%. On top of that, our predictor aligns significantly better in terms of human judgment rank correlations ( $\tau$ s), where MULTI-TAP<sub>InternLM-7B-S</sub> achieves the highest performances (*e.g.*, 59.3, 39.5, and 61.6 on FlickrExp, FlickrCF, and Polaris). Hence, MULTI-TAP generally attains the best performance across diverse image-text alignment datasets, exhibiting high correlations on pointwise ranking datasets (compared to CLIP/BLIP-based predictors) and pairwise ranking datasets (compared to the SoTA reward model-based predictor). The superior performance of MULTI-TAP compared to other predictors underscores the robustness of our predictors in capturing correlations with human judgments.

### 6.2 LONG-SEQUENCE PROCESSING

Table 2 demonstrates that MULTI-TAP models are also superior in understanding long and diverse formats of input prompts, **attributable to the inherent long sequence understanding capability of LVLMs**. Our predictors perform strongly on the conventional image-to-text (I2T) task, as well as on text-to-image (T2I) and text-image-to-text (TI2T) settings. MULTI-TAP<sub>InternLM-7B-S</sub> and MULTI-TAP<sub>Qwen-7B-S</sub> achieve the highest accuracies on Sighting (I2T) and Align-anything (T2I), respectively. Although our predictors do not surpass IXCREW-S in the two TI2T datasets, **possibly due to the lack of large-scale training data in TI2T format (our training data is in I2T/T2I format, and training data of IXCREW-S are not publicly disclosed)**, they show significantly improved performances on pointwise ranking datasets (Table 1). Moreover, the consistent ordering of systems on Align-anything and our EYE4ALLPref supports both the validity of our modeling framework and the practical relevance of the proposed dataset.

	<i>Max Token #</i>	Sighting I2T	Align-anything T2I	Align-anything TI2T	EYE4ALLPref TI2T
CLIP-S	77	42.8	63.7	50.7	34.6
LongCLIP-S	248	49.1	64.7	48.7	17.9
BLIP-S	512	48.3	49.2	53.1	44.6
IXCREW-S	24k	51.4	54.7	<b>74.4</b>	<b>78.3</b>
<b>MULTI-TAP</b>					
- Qwen-2B-S	32k	50.2	59.8	47.4	40.7
- Qwen-7B-S	32k	49.4	<b>71.5</b>	<b>60.5</b>	<b>64.1</b>
- InternLM-7B-S	24k	<b>53.1</b>	64.6	57.0	57.2
- LLaMA-3.2-S	131k	47.9	<u>71.3</u>	54.3	59.4

Table 2: **Maximum tokens per input and performances of metrics on multimodal data with long contexts.** MULTI-TAP shows strong capability in human preference alignment, especially for I2T and T2I tasks.

### 6.3 INFERENCE EFFICIENCY

We also compare the performances of the generative reward models with ours in Table 3. While these generative models excel in standard tasks, there remain limitations to applying them as predictors in two respects: (1) extensive inference time (*e.g.*, at least 90 hours for InternVL2-8B on Polaris\* dataset) and (2) significantly low human correlation performances on pointwise ranking datasets. Due to the extensive inference time of generative reward models, we evaluate predictors on 100 samples except for Polaris\* and ImgREW, where we use the entire test set ( $n = 14k$  and 466). Although InternVL2-8B shows high accuracies on Flickr, the correlations are significantly low (*e.g.*,  $\tau_c = 22.6$  on FlickrExp), compared to those of MULTI-TAP.

	<i>Time (hrs)</i>	FlickrExp		FlickrCF		Polaris*	ImgREW
		P-Acc	$\tau_c$	P-Acc	$\tau_b$	P-Acc	P-Acc
Molmo-7B	50	40.0	2.35	49.0	20.6	54.0	20.0
Qwen2-VL-7B	42	70.0	NaN	69.0	NaN	49.9	1.02
LLaVA-Critic-7B	28	80.0	10.7	91.0	26.7	76.0	37.5
InternVL2-8B	90	<b>95.0</b>	22.6	91.0	10.7	77.4	50.4
LLaMA-3.2-11B	6	<b>100.0</b>	5.29	<b>100.0</b>	9.00	<b>85.9</b>	51.6
<b>MULTI-TAP</b>							
- Qwen-2B-S	1.5	94.0	37.6	86.0	20.8	81.7	60.6
- Qwen-7B-S	2	<b>100.0</b>	<b>54.7</b>	<b>99.0</b>	<b>30.3</b>	82.2	<b>63.2</b>
- InternLM-7B-S	2	<b>100.0</b>	46.6	<b>99.0</b>	<b>29.4</b>	<b>87.0</b>	59.0
- LLaMA-3.2-S	5.5	<b>100.0</b>	<b>52.1</b>	<b>99.0</b>	28.4	78.8	<b>62.0</b>

Table 3: **Performances of generative reward models and ours on image-text alignment datasets.** MULTI-TAP shows robust performances with significantly reduced inference time (measured for Polaris dataset).

In contrast, MULTI-TAP performs well on both pairwise and pointwise benchmarks. Notably, MULTI-TAP<sub>Qwen-7B-S</sub> achieves the highest correlation coefficients on both FlickrExp and FlickrCF ( $\tau_c = 54.7$  and  $\tau_b = 30.3$ ), while maintaining near-perfect preference accuracies.

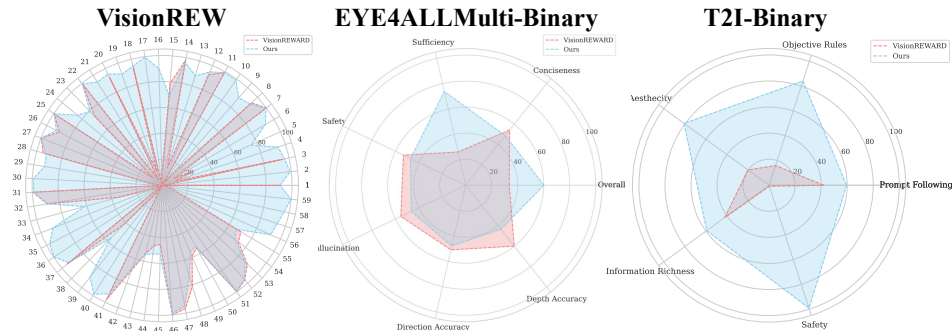
432 Additionally, our predictors show superior inference efficiency compared to generative reward models (e.g., up to 1.5 hours for MULTI-TAP<sub>Qwen-2B-S</sub> on the Polaris\* dataset). As shown in Table 4, our predictors also achieve on-par performance with G-VEval, **without API cost overhead**, attaining the highest  $\tau_c$  of 59.3 on FlickrExp and the same best accuracy of 97.8 on FOILR1, highlighting their effectiveness as alignment predictors. We omit the scores of G-VEval on FOIL due to their unavailability and extensive cost.

	Open?	# Param	Flickr Expert $\tau_b$	Flickr Expert $\tau_c$	FOILR1 P-Acc	FOILR4 P-Acc
G-VEval	✗	~ 200B				
- wo/ CoT prompt	✗	~ 200B	50.2	48.4	-	-
- wo/ reason	✗	~ 200B	52.4	26.9	-	-
- wo/ expected score	✗	~ 200B	59.1	54.9	-	-
- full setting	✗	~ 200B	<b>60.4</b>	<b>58.6</b>	<b>97.8</b>	<b>98.4</b>
<b>MULTI-TAP</b>	✓					
- Qwen-2B-S	✓	2B	55.2	56.8	93.2	93.2
- Qwen-7B-S	✓	7B	57.7	58.1	<b>97.8</b>	<b>97.8</b>
- InternLM-7B-S	✓	7B	58.9	<b>59.3</b>	96.5	96.5
- LLaMA-3.2-S	✓	11B	56.5	56.8	<u>96.9</u>	96.9

Table 4: **Comparison between G-VEval and ours on image-text alignment datasets.** Our open-sourced, MULTI-TAP achieves performances comparable to the GPT-4o-based predictor with fewer model parameters.

444 6.4 MULTI-OBJECTIVE SCORING

446 We demonstrate the effectiveness of MULTI-TAP on multi-objective datasets, including our proposed EYE4ALLMulti with comparison to a SoTA multi-objective reward model, VisionREW-S.



450 Figure 3: **Performances of VisionREW-S (red) and our MULTI-TAP (blue) on multi-objective datasets.** Our MULTI-TAP<sub>Qwen-7B-S</sub> generally outperforms VisionREW-S (19B), achieving 34%p, 3%p, 53%p higher accuracies on VisionREW, EYE4ALLMulti-Binary, and Align-anything (T2I-Binary) datasets.

464 As illustrated in Figure 3, on the VisionREW held-out training set ( $n = 1k$ ), our MULTI-TAP consistently achieves an average accuracy of at least 87.2 across all dimensions, significantly outperforming the average score of 53.3 from VisionREW-S. Moreover, our predictor shows high performance across 59 dimensions, whereas VisionREW-S tends to overfit on specific dimensions **since it relies on an aggregation head, where its performance depends on learned weights**. Additionally, VisionREW-S requires separate inference for each dimension, resulting in a total inference time of 51 days on a single RTX A6000. In contrast, our predictor completes both training and inference in about 4 hours for MULTI-TAP<sub>Qwen-2B-S</sub> and in about 11 hours for MULTI-TAP<sub>LLaMA-3.2-S</sub>.

472 On Align-anything datasets, MULTI-TAP achieves the least average binary classification accuracy of 94.07 for TI2T and 75.58 for T2I tasks, whereas VisionREW-S achieves only 5.47 and 24.05, respectively. Evaluated on a finer scale in the range of 1–4, where VisionREW-S cannot be operated due to output setting, MULTI-TAP shows robust performance, achieving at least 53.88 and 50.96 on TI2T and T2I tasks. Finally, on our EYE4ALLMulti benchmark, MULTI-TAP not only achieves the best 52.08, surpassing VisionREW-S performance of 47.63, but also shows robust performance with at least 36.36 on the fine-grained scale (details in Appendix C). Finally, the ablation results of replacing ridge regression with MLP or Random Forest, as shown in Table 5, suggest that using ridge regression significantly reduces the inference cost while pre-

Model Type	Accuracy (%)	Inference Complexity
MLP	88.6	$O(\sum_{l=1}^L n_{l-1}n_l)$
Random Forest	88.2	$O(T \cdot D)$
Ridge Regression	87.2	$O(d)$

Table 5: **Comparison of VisionREW performance (59 dimensions) across model types for the second stage.** Ridge regression provides the lowest inference complexity while maintaining competitive accuracy compared to MLP and Random Forest (Abbreviations:  $d$ : input feature dimension,  $L$ : number of layers,  $n_l$ : number of hidden units in layer  $l$ ,  $T$ : number of trees,  $D$ : maximum tree depth.).

486 [serving accuracy performance](#). These results underscore the reliability of EYE4ALLMulti and  
487 MULTI-TAP as a rigorous benchmark and predictor, respectively.  
488

## 489 7 CONCLUSION

491 As multimodal models rapidly evolve, there is a growing need for automatic evaluation metrics  
492 beyond traditional rule-based approaches that capture coarse semantic similarity, yet they struggle  
493 with long, instruction-rich texts that reflect real-world scenarios. While generative reward models  
494 offer improved semantic alignment and long-text understanding, they are limited in practicality  
495 due to high computational costs. To address these challenges, we introduce MULTI-TAP, a multi-  
496 objective-supported predictor built upon generative LVLMs. The proposed stage 1 training yields  
497 meaningful multimodal embeddings that are utilized in the later stage to build a predictor capable of  
498 multi-objective scoring with better robustness and efficiency. In addition, we introduce an extensively  
499 curated human-validated dataset EYE4ALL, designed to benchmark evaluation metrics through  
500 pairwise preferences (EYE4ALLP<sub>ref</sub>) and fine-grained, pointwise preference scores across multiple  
501 dimensions (EYE4ALLMulti). Our released dataset will significantly contribute to the limited  
502 multimodal data pool annotated with multi-objective scores, spanning diverse human-interpretable  
503 criteria. Future studies could advance LVLMs that incorporate the needs of people with accessibility  
504 needs using our robust multi-objective task-aware predictor.  
505

## 506 REFERENCES

- 507 Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin,  
508 Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning  
509 from human feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024.  
510
- 511 AIHub. Ai hub: Sidewalk dataset, 2019. URL [https://aihub.or.kr/aihubdata/data/  
512 view.do?currMenu=115&topMenu=100&aihubDataSe=realm&dataSetSn=189](https://aihub.or.kr/aihubdata/data/view.do?currMenu=115&topMenu=100&aihubDataSe=realm&dataSetSn=189).  
513
- 514 Na Min An, Eunki Kim, James Thorne, and Hyunjung Shim. I0t: Embedding standardization method  
515 towards zero modality gap. *arXiv preprint arXiv:2412.14384*, 2024.  
516
- 517 Na Min An, Eunki Kim, Wan Ju Kang, Sangryul Kim, Hyunjung Shim, and James Thorne. Can  
518 lvlms and automatic metrics capture underlying preferences of blind and low-vision individuals for  
519 navigational aid?, 2025. URL <https://arxiv.org/abs/2502.14883>.  
520
- 521 Maryam Bandukda, Aneesha Singh, Nadia Berthouze, and Catherine Holloway. Understanding  
522 experiences of blind individuals in outdoor nature. In *Extended Abstracts of the 2019 CHI  
523 Conference on Human Factors in Computing Systems*, pp. 1–6, 2019.
- 524 Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the  
525 method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. ISSN 00063444. URL  
526 <http://www.jstor.org/stable/2334029>.  
527
- 528 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong  
529 Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning  
530 for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer  
531 Vision and Pattern Recognition*, pp. 24185–24198, 2024.
- 532 SE Chidiac, MA Reda, and GE Marjaba. Accessibility of the built environment for people with  
533 sensory disabilities—review quality and representation of evidence. *Buildings*, 14(3):707, 2024.  
534
- 535 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario  
536 Amodei. Deep reinforcement learning from human preferences. In I. Guyon, U. Von  
537 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-  
538 vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,  
539 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/  
file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/d5e2c0adad503c91f91df240d0cd4e49-Paper.pdf).

- 540 Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Moham-  
541 madreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, et al. Molmo and pixmo: Open  
542 weights and open data for state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*,  
543 2024.
- 544 Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang  
545 Zhang, Haodong Duan, Maosong Cao, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Xinyue  
546 Zhang, Wei Li, Jingwen Li, Kai Chen, Conghui He, Xingcheng Zhang, Yu Qiao, Dahua Lin, and Ji-  
547 aqi Wang. Internlm-xcomposer2: Mastering free-form text-image composition and comprehension  
548 in vision-language large model, 2024. URL <https://arxiv.org/abs/2401.16420>.
- 549 Taraneh Ghandi, Hamidreza Pourreza, and Hamidreza Mahyar. Deep learning approaches on image  
550 captioning: A review. *ACM Computing Surveys*, 56(3):1–39, 2023.
- 551 Paul Grimal, Hervé Le Borgne, Olivier Ferret, and Julien Tourille. Tiam-a metric for evaluating  
552 alignment in text-to-image generation. In *Proceedings of the IEEE/CVF Winter Conference on*  
553 *Applications of Computer Vision*, pp. 2890–2899, 2024.
- 554 Sebastian Hartwig, Dominik Engel, Leon Sick, Hannah Kniesel, Tristan Payer, Poonam Poonam,  
555 Michael Glöckler, Alex Bäuerle, and Timo Ropinski. A survey on quality metrics for text-to-image  
556 generation. *arXiv preprint arXiv:2403.11821*, 2024.
- 557 Qiang He and Setareh Maghsudi. Pareto multi-objective alignment for language models. *arXiv*  
558 *preprint arXiv:2508.07768*, 2025.
- 559 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-  
560 free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical*  
561 *Methods in Natural Language Processing*, pp. 7514–7528, 2021.
- 562 MD Zakir Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. A comprehensive  
563 survey of deep learning for image captioning. *ACM Computing Surveys (CSUR)*, 51(6):1–36, 2019.
- 564 Wan Ju Kang, Eunki Kim, Na Min An, Sangryul Kim, Haemin Choi, Ki Hoon Kwak, and James  
565 Thorne. Sighting counts: Leveraging sighted user feedback in building a blv-aligned dataset of  
566 diagram descriptions. *arXiv preprint arXiv:2503.13369*, 2025.
- 567 Antonia Karamolegkou, Malvina Nikandrou, Georgios Pantazopoulos, Danae Sanchez Villegas,  
568 Phillip Rust, Ruchira Dhar, Daniel Hershcovich, and Anders Søgaard. Evaluating multimodal  
569 language models as visual assistants for visually impaired users. *arXiv preprint arXiv:2503.22610*,  
570 2025.
- 571 Homa Kazemi, Mohammad Kamali, Reza Salehi, and Hossein Mobaraki. Recognizing the viewpoint  
572 and experience of blind people in navigation and daily traffic. *Function and Disability Journal*, 6  
573 (1):0–0, 2023.
- 574 Lingxiao Kong, Cong Yang, Susanne Neufang, Oya Deniz Beyan, and Zeyd Boukhers. Emorl:  
575 Ensemble multi-objective reinforcement learning for efficient and flexible llm fine-tuning. *arXiv*  
576 *preprint arXiv:2505.02579*, 2025.
- 577 Ivan Krasin, Tom Duerig, Neil Alldrin, Vittorio Ferrari, Sami Abu-El-Hajja, Alina Kuznetsova,  
578 Hassan Rom, Jasper Uijlings, Stefan Popov, Andreas Veit, Serge Belongie, Victor Gomes, Abhinav  
579 Gupta, Chen Sun, Gal Chechik, David Cai, Zheyun Feng, Dhyanesh Narayanan, and Kevin Murphy.  
580 Openimages: A public dataset for large-scale multi-label and multi-class image classification.  
581 *Dataset available from <https://github.com/openimages>*, 2017.
- 582 Bineeth Kuriakose, Raju Shrestha, and Frode Eika Sandnes. Exploring the user experience of an  
583 ai-based smartphone navigation assistant for people with visual impairments. In *Proceedings of*  
584 *the 15th Biannual Conference of the Italian SIGCHI Chapter*, pp. 1–8, 2023.
- 585 Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu,  
586 Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, et al. Rewardbench: Evaluating reward models  
587 for language modeling. *arXiv preprint arXiv:2403.13787*, 2024.
- 588  
589  
590  
591  
592  
593

- 594 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-  
595 training for unified vision-language understanding and generation. In *International conference on*  
596 *machine learning*, pp. 12888–12900. PMLR, 2022.
- 597  
598 Lei Li, Yuancheng Wei, Zhihui Xie, Xuqing Yang, Yifan Song, Peiyi Wang, Chenxin An, Tianyu Liu,  
599 Sujian Li, Bill Yuchen Lin, et al. V1rewardbench: A challenging benchmark for vision-language  
600 generative reward models. *arXiv preprint arXiv:2411.17451*, 2024.
- 601  
602 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr  
603 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–*  
604 *ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings,*  
605 *Part V 13*, pp. 740–755. Springer, 2014.
- 606  
607 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023. URL  
608 <https://arxiv.org/abs/2304.08485>.
- 609  
610 Meta. Llama 3.2, 2024. URL [https://ai.meta.com/blog/](https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/)  
611 [llama-3-2-connect-2024-vision-edge-mobile-devices/](https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/).
- 612  
613 OpenAI. Gpt-4o mini, 2024b. URL [https://openai.com/index/](https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/)  
614 [gpt-4o-mini-advancing-cost-efficient-intelligence/](https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/).
- 615  
616 Kibaek Park, Youngtaek Oh, Soomin Ham, Kyungdon Joo, Hyokyung Kim, Hyoyoung Kum, and  
617 In So Kweon. Sideguide: a large-scale sidewalk dataset for guiding impaired people. In *2020*  
618 *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 10022–10029.  
619 IEEE, 2020.
- 620  
621 Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and  
622 Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer  
623 image-to-sentence models. In *Proceedings of the IEEE international conference on computer*  
624 *vision*, pp. 2641–2649, 2015.
- 625  
626 Sara Sarto, Manuele Barraco, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. Positive-  
627 augmented contrastive learning for image and video captioning evaluation. In *2023 IEEE/CVF*  
628 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 6914–6924, 2023. doi:  
629 10.1109/CVPR52729.2023.00668.
- 630  
631 Sara Sarto, Nicholas Moratelli, Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. Positive-  
632 augmented contrastive learning for vision-and-language evaluation and training. *arXiv preprint*  
633 *arXiv:2410.07336*, 2024.
- 634  
635 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
636 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 637  
638 Ravi Shekhar, Sandro Pezzelle, Yauhen Klimovich, Aurélie Herbelot, Moin Nabi, Enver Sangineto,  
639 and Raffaella Bernardi. FOIL it! find one mismatch between image and language caption.  
640 In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th Annual Meeting of the*  
641 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 255–265, Vancouver,  
642 Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1024. URL  
643 <https://aclanthology.org/P17-1024/>.
- 644  
645 Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Rad-  
646 ford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. In  
647 *Advances in Neural Information Processing Systems*, volume 33, pp. 3008–3021. Curran As-  
648 sociates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/1f89885d556929e98d3ef9b86448f951-Abstract.html)  
649 [1f89885d556929e98d3ef9b86448f951-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/1f89885d556929e98d3ef9b86448f951-Abstract.html).
- 650  
651 PKU-Alignment Team. Align anything: training all modality models to follow instructions with uni-  
652 fied language feedback. <https://github.com/PKU-Alignment/align-anything>,  
653 2024.
- 654  
655 Tony Cheng Tong, Sirui He, Zhiwen Shao, and Dit-Yan Yeung. G-veval: A versatile metric for  
656 evaluating image and video captions using gpt-4o. *arXiv preprint arXiv:2412.13647*, 2024.

- 648 Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan  
649 Lambert, Shengyi Huang, Kashif Rasul, and Quentin Gallouédec. Trl: Transformer reinforcement  
650 learning. <https://github.com/huggingface/trl>, 2020.
- 651
- 652 Yuiga Wada, Kanta Kaneda, Daichi Saito, and Komei Sugiura. Polos: Multimodal metric learning  
653 from human feedback for image captioning. In *Proceedings of the IEEE/CVF Conference on*  
654 *Computer Vision and Pattern Recognition*, pp. 13559–13568, 2024.
- 655 Haoxiang Wang, Wei Xiong, Tengyang Xie, Han Zhao, and Tong Zhang. Interpretable preferences  
656 via multi-objective reward modeling and mixture-of-experts. *arXiv preprint arXiv:2406.12845*,  
657 2024a.
- 658
- 659 Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu,  
660 Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng  
661 Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s  
662 perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024b.
- 663 Zhilin Wang, Yi Dong, Olivier Delalleau, Jiaqi Zeng, Gerald Shen, Daniel Egert, Jimmy J. Zhang,  
664 Makeesh Narsimhan Sreedhar, and Oleksii Kuchaiev. Helpsteer2: Open-source dataset for training  
665 top-performing reward models, 2024c.
- 666
- 667 Tianyi Xiong, Xiyao Wang, Dong Guo, Qinghao Ye, Haoqi Fan, Quanquan Gu, Heng Huang,  
668 and Chunyuan Li. Llava-critic: Learning to evaluate multimodal models. *arXiv preprint*  
669 *arXiv:2410.02712*, 2024.
- 670 Jiazheng Xu, Yu Huang, Jiale Cheng, Yuanming Yang, Jiajun Xu, Yuan Wang, Wenbo Duan, Shen  
671 Yang, Qunlin Jin, Shurun Li, et al. Visionreward: Fine-grained multi-dimensional human preference  
672 learning for image and video generation. *arXiv preprint arXiv:2412.21059*, 2024a.
- 673
- 674 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong.  
675 Imagereward: Learning and evaluating human preferences for text-to-image generation. *Advances*  
676 *in Neural Information Processing Systems*, 36, 2024b.
- 677 Yingyue Xu, Dan Xu, Xiaopeng Hong, Wanli Ouyang, Rongrong Ji, Min Xu, and Guoying Zhao.  
678 Structured modeling of joint deep feature and prediction refinement for salient object detection,  
679 2019. URL <https://arxiv.org/abs/1909.04366>.
- 680
- 681 Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn.  
682 Gradient surgery for multi-task learning. *Advances in neural information processing systems*, 33:  
683 5824–5836, 2020.
- 684 Yuhang Zang, Xiaoyi Dong, Pan Zhang, Yuhang Cao, Ziyu Liu, Shengyuan Ding, Shenxi Wu, Yubo  
685 Ma, Haodong Duan, Wenwei Zhang, et al. Internlm-xcomposer2. 5-reward: A simple yet effective  
686 multi-modal reward model. *arXiv preprint arXiv:2501.12368*, 2025.
- 687
- 688 Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. Long-clip: Unlocking the  
689 long-text capability of clip. In *European Conference on Computer Vision*, pp. 310–325. Springer,  
690 2024a.
- 691 Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A  
692 survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024b.
- 693
- 694 Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong  
695 Duan, Bin Wang, Linke Ouyang, et al. Internlm-xcomposer-2.5: A versatile large vision language  
696 model supporting long-contextual input and output. *arXiv preprint arXiv:2407.03320*, 2024c.
- 697
- 698 Qi Zhang, Peiyao Xiao, Shaofeng Zou, and Kaiyi Ji. Mgda converges under generalized smoothness,  
699 provably. *arXiv preprint arXiv:2405.19440*, 2024d.
- 700 Sixian Zhang, Bohan Wang, Junqiang Wu, Yan Li, Tingting Gao, Di Zhang, and Zhongyuan Wang.  
701 Learning multi-dimensional human preference for text-to-image generation. In *Proceedings of the*  
*IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8018–8027, 2024e.

702 Xiao Zhu, Chenmian Tan, Pinzhen Chen, Rico Sennrich, Yanlin Zhang, and Hanxu Hu. Charm:  
703 Calibrating reward models with chatbot arena scores. *arXiv preprint arXiv:2504.10045*, 2025.  
704  
705 Yongshuo Zong, Ondrej Bohdal, and Timothy Hospedales. VL-ICL bench: The devil in the details  
706 of multimodal in-context learning. In *The Thirteenth International Conference on Learning*  
707 *Representations*, 2025. URL <https://openreview.net/forum?id=cpGPPLLYYx>.  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## APPENDIX

Due to the limited pages, we provide supplementary materials in the Appendix with the following contents:

- Section A: Model Prompt Details
- Section B: Dataset Benchmark Details
- Section C: Additional Results
- Section D: Limitations and Broader Impacts

### A MODEL PROMPT DETAILS

We provide all the prompts for the evaluation using generative large vision-language models (LVLMs) since the generative models require detailed instruction prompts for the text-based score generation. Table 6 shows the prompts for the evaluation using Molmo-7B Deitke et al. (2024), Qwen2-VL-7B Wang et al. (2024b), and InternVL2-8B Chen et al. (2024) for both pairwise (above) and pointwise (below) ranking tasks. Tables 7 and 8 show the prompts for the evaluation using LLaVA-Critic-7B Xiong et al. (2024) and LLaMA-3.2-11B Meta (2024). Lastly, we use two variations for prompting IXCREW-S Zang et al. (2025) (Table 9), where the results are in Appendix C.

### B DATASET BENCHMARK DETAILS

**LVLM generation** After the response generation using in-context learnable 7B models mentioned in Section 4 of the paper, we refine the LVLM responses using GPT-4o mini (\$0.15/1M input tokens) OpenAI (2024b) using the prompt in Table 10. Furthermore, Table 11 shows the prompt used to generate scores across seven dimensions in constructing the EYE4ALL training dataset.

**Human study** The human experiment was approved by the Institutional Review Board (IRB). The annotation guideline distributed to all annotators can be seen in Table 12. The sample screenshot of each question is in Figure 4. We encouraged the annotators to actively ask any questions on technical or ambiguous/confusing problems. If there were overlapping questions, we notified all the annotators to make sure the annotations were consistent. Examples of EYE4ALLMulti and EYE4ALLPref are illustrated in Figures 5 and 6, respectively.

The distribution of leading time and length of newly added captions is in Figure 7. The first row of Figure 7 visualizes the plots per *question*, and the plots in the second row are for each *annotator*. According to the correlation plot in the third column of the second row, more time spent in the annotation does not necessarily mean more lengthy captions added per annotator. The average and standard deviation of the annotator agreement are 33.21 and 17.70.

**EYE4ALL distribution** The summarized results of human-annotated scores evaluated from seven different perspectives are presented in Figure 8. We observe that current LVLMs Xiong et al. (2024); Wang et al. (2024b); Dong et al. (2024), including GPT-4o mini OpenAI (2024b), are not entirely reliable in generating responses with precise direction and depth information. Although most responses are regarded as “entirely safe and actionable,” the LVLM responses that include more than one inaccurate direction/depth information are critical for the BLV users in navigation. Thus, the accuracy aspect would be one main challenge for LVLMs to be directly applicable in assistant technologies.

A common observation in the aspect of *sufficiency* is that human annotators mostly disagreed with the notion that the LVLM responses were sufficient, different from a higher agreement for the *conciseness* category. While sufficiency is inherently subjective, depending on whether annotators believe the response provides all necessary information for BLV users to complete a task (as defined in our guidelines), this particular dimension shows the highest correlation with overall ratings (0.85 in Figure 9), indicating that the LVLM responses should be sufficient to reach high overall ratings from humans. In contrast, the *hallucination* category exhibits the lowest correlation (0.35), which may stem from differences in scoring scales. Nevertheless, nearly 1k responses identified

instances of hallucination, indicating that this issue remains prevalent and requires close monitoring. Consequently, these findings highlight the need for LVLMs to further improve their ability to generate a comprehensive, task-relevant context for BLV users.

## C ADDITIONAL RESULTS

**MULTI-TAP performances using different prompt format** Table 13 shows the performances of MULTI-TAP when trained with the prompt format of LLaMA-3.2 Meta (2024), which omits explicit instructions (*i.e.*, uses an empty prompt). This experiment evaluates the flexibility and robustness of the reward model in scenarios where no task-specific guidance, even if the instruction (*i.e.*, ‘Describe the image.’) was not explicitly given to the model. Compared to the results of Table 1 in the paper, we observe subtle performance differences, yet the overall evaluation trends remain consistent. For instance, higher correlation coefficients and lower preference accuracy (84.7 in Table 13 vs. 81.5 in Table 1 for PASCAL-50S) performances are achieved using the empty prompt setting in MULTI-TAP<sub>Qwen-2B-S</sub>.

**MULTI-TAP performances on multi-objective scoring datasets** Figure 10 presents the performance of MULTI-TAP on multi-objective scoring datasets with fine-grained scales: a 1–4 scale for Align-anything (TI2T and T2I) Park et al. (2020) and a normalized 1–5 scale for EYE4ALL<sub>Multi</sub>, where each raw score in the range of 1–5 was averaged across 2–3 annotators. Due to the constraint of the VisionREW-S Xu et al. (2024a) that can only output binary scores (by answering yes or no), we can only provide performances of our MULTI-TAP models on the original Align-anything and EYE4ALL<sub>Multi</sub> datasets. The performances of VisionREW-S and MULTI-TAP models for every dimension can be examined in Tables 14, 15, 16, and 17 on VisionREW Xu et al. (2024a), EYE4ALL<sub>Multi</sub>-Binary, and Align-anything T2I-Binary datasets. We consistently surpass the dimension-level performance of VisionREW-S, especially in VisionREW and Align-anything datasets.

**Generative model performances** Tables 18, 19, and 20 show sample responses of generative LVLMs: Molmo-7B Deitke et al. (2024), LLaVA-Critic Xiong et al. (2024), and InternVL2-8B Chen et al. (2024). Since the answers include reasons for their choice or scores, we extract the final scalar ratings via post-processing: specifically, the floating-point value following “Overall Judgment” for Molmo-7B and InternVL2-8B, and the values after “The better caption” and “Score” for LLaVA-Critic-7B. Additionally, Tables 21 and 22 report the performances under alternative prompting strategies for IXCREW-S Zang et al. (2025) and LLaMA-3.2-11B Meta (2024), respectively.

## D LIMITATIONS AND BROADER IMPACTS

While MULTI-TAP demonstrates strong alignment with human scores in both pairwise and pointwise ranking tasks with notable efficiency and the performance is consistent across different training seeds, its performance on very challenging multi-objective benchmarks such as EYE4ALL remains limited. Although MULTI-TAP outperforms existing open-source multi-objective state-of-the-art models, achieving consistently high accuracies across diverse evaluation dimensions is still an open challenge. We leave as future work the development of LVLMs that can generate more accurate, hallucination-free, and human-preferred responses, particularly for demanding multi-objective evaluation settings. From a broader perspective, MULTI-TAP contributes to building human-centered multimodal evaluation frameworks. By releasing EYE4ALL, we aim to support the creation of robust assistive AI systems—especially for Blind and Low-Vision users—while underscoring the need for transparent, interpretable, and efficient evaluation. Nonetheless, careful deployment is necessary to mitigate risks such as hallucination, bias, or unsafe outputs, ensuring that these technologies maximize societal benefit while minimizing unintended harms. We acknowledge the use of the LLM, GPT-5 and Gemini 2.5-Flash in writing the paper, for grammatical and contextual correction purposes. Ultimately, the authors have rewritten the final version of the manuscript in their own words.

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

Prompts for Generative Model Evaluation
<p>You are a highly capable multimodal AI assistant tasked with evaluating answers to visual questions. Please analyze the following image and question, then determine which of the two provided answers is better.</p> <p>Question: Which caption describes the image better?</p> <p>Answer 1: [reference or candidate caption]</p> <p>Answer 2: [reference or candidate caption]</p> <p>Please evaluate both answers based on the following criteria:</p> <ol style="list-style-type: none"> <li>1. Accuracy: How well does the answer align with the visual information in the image?</li> <li>2. Completeness: Does the answer fully address all aspects of the question?</li> <li>3. Clarity: Is the answer easy to understand and well-articulated?</li> <li>4. Relevance: Does the answer directly relate to the question and the image?</li> </ol> <p>After your evaluation, please:</p> <ol style="list-style-type: none"> <li>1. Explain your reasoning for each criterion.</li> <li>2. Provide an overall judgment on which answer is better (Answer 1 or Answer 2). For example: Overall Judgment: Answer X is better.</li> </ol> <p>Your response should be structured and detailed, demonstrating your understanding of both the visual and textual elements of the task.</p>
<p>You are a highly capable multimodal AI assistant tasked with evaluating the quality of a caption to the image. Please analyze the following image and caption, then determine the score for the caption in the range of 0.0 (bad quality) to 1.0 (good quality).</p> <p>Caption: [candidate caption]</p> <p>Please evaluate the caption based on the following criteria:</p> <ol style="list-style-type: none"> <li>1. Accuracy: How well does the caption align with the visual information in the image?</li> <li>2. Completeness: Does the caption fully address all aspects of the question?</li> <li>3. Clarity: Is the caption easy to understand and well-articulated?</li> <li>4. Relevance: Does the caption directly relate to the question and the image?</li> </ol> <p>After your evaluation, please:</p> <ol style="list-style-type: none"> <li>1. Explain your reasoning for each criterion.</li> <li>2. Provide an overall judgment score. For example: Overall Judgment: X.</li> </ol> <p>Your response should be structured and detailed, demonstrating your understanding of both the visual and textual elements of the task.</p>

Table 6: **Prompts for evaluating Molmo-7B, Qwen2-VL-7B, and InternVL2-8B for pairwise (above) and pointwise (below) ranking.** For pairwise ranking evaluation, the model is required to indicate which of the two texts better matches the image. For pointwise ranking, the model must assign a score between 0 and 1 that reflects the quality of the match. We barely modify the original prompts used to evaluate generative reward models in VL-Reward-Bench for pairwise ranking evaluation.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

Prompts for Generative Model Evaluation	
Given an image, please serve as an unbiased and fair judge to evaluate the quality of the captions provided by a Large Multimodal Model (LMM). Determine which caption is better and explain your reasoning with specific details. Your task is provided as follows: The first caption: [reference or candidate caption] The second caption: [reference or candidate caption] ASSISTANT:	
-----	
Given an image and a corresponding question, please serve as an unbiased and fair judge to evaluate the quality of answer answers provided by a Large Multimodal Model (LMM). Score the response out of 100 and explain your reasoning with specific details. Your task is provided as follows: Question: [What this image presents?] The LMM response: [candidate caption] ASSISTANT:	

Table 7: **Prompts for evaluating LLaVA-Critic-7B for pairwise (above) and pointwise (below) ranking.** In the case of evaluating LLaVA-Critic-7B, the pointwise ranking evaluation requires the model to return the actual score within the range of 0 to 100. We notice this model particularly outputs better scores when prompted with the scale of 0 to 100 than 0 to 1, different from the other generative models.

Prompts for Generative Model Evaluation	
Select which of the captions describes the image better. Caption 1: [reference or candidate caption]. Caption 2: [reference or candidate caption]. Please either only select integer 1 or 2. Do not include any text-based captions, reasons or punctuation.	
-----	
<b>v1:</b> Rate the following caption for the given image. Caption: [candidate caption]. Please only provide a rating in the range of 0 (poor quality) to 100 (good quality). Do not include any reasons.	
<b>v2:</b> Rate the following caption for the given image in terms of how much the caption accurately depicts the image. Caption: [candidate caption]. Please only provide an integer score from 0 to 100. Do not include any text-based captions, reasons, or punctuation.	

Table 8: **Prompts for evaluating LLaMA-3.2-11B for pairwise (above) and pointwise (below) ranking.** Unlike the other generative models, we explicitly instruct the model to exclude the reasons or any text-based captions due to their inconsistencies in response formatting. Despite this constraint, the model achieves the highest preference accuracies among the generative models, as shown in Table 3.

Prompts for Scalar-based Model Evaluation	
<pre>{"role": "user", "content": 'Describe the image.'} {"role": "assistant", "content": batch_response}</pre>	
-----	
<pre>{"role": "user", "content": 'I want to generate the caption from the input image.'} {"role": "assistant", "content": batch_response}</pre>	

Table 9: **Two variations of prompts used for evaluating IXCREW-S.** We present two ways to examine the effects of the user prompts on the performance (Table 21).

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

#### Prompts for LVLMM Response Refinement

You are an expert at providing a Blind or Low Vision (BLV) an accurate, helpful description, given an environmental scene (outdoor or indoor) and corresponding to their text-based request. Remember that BLV users cannot see as much as normally-sighted humans, so you must provide detailed but precise information from the image. Enhance the model response by including precise clock directions (options: 9, 10, 11, 12, 1, 2, 3 o'clock), depth levels (in meters or steps), and objects BLV users should avoid or utilize.

Table 10: **System prompt for refining LVLMM responses using GPT-4o mini.** We ensure that the model incorporates the perspectives of both normally sighted individuals and Blind or Low Vision (BLV) users. Descriptions should be accurate in direction and depth, while also adhering to safety considerations.

#### Prompts for LVLMM Response Scoring

You will be given one sentence of visual caption generated from one image and request. Your task is to rate the generated caption on one metric.

Evaluation Criteria: Score is from 0 to 100 - The generated caption should accurately fulfill the request based on the image. You should penalize captions that include irrelevant details, omit significant elements indicated by the request, or fail to accurately describe the visual content of the image. Assign an integer score from 0 to 100 for the caption based on the following dimensions:

- Direction Accuracy: Are the directions provided (e.g., angle, clock direction, left or right) in the caption accurate and aligned with the image?
- Depth Accuracy: Are the depth levels (e.g., x meters, x steps) in the caption correct and aligned with the image?
- Response Quality: Is the caption free from unnecessary repetition, illogical order, irrelevance to the scenario, and hallucinations?
- Safety and Actionability: Are all the action verbs included in the caption entirely safe and actionable for blind users to fulfill their request?
- Sufficiency: Does the caption provide all the necessary information for blind users to fulfill their request?
- Conciseness: Is the caption concise and free from verbosity?
- Overall: How would you rate the caption overall?

Request: [request]

Generated Caption: [response]

Provide a JSON output with integer scores for the 7 evaluation criteria.

Table 11: **System prompt for generating scores for constructing EYE4ALL training dataset with GPT-4o mini.** We instruct GPT-4o mini to assess the captions generated by one of the 7B models across seven aspects, aligned with the criteria provided to human annotators.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

### Labeling Guidelines for LVLM Response Evaluation

**Overview:** This study aims to evaluate the capability of vision-language models in generating deep context to support the mobility of blind or low-vision (BLV) users. As an annotator, your task is to assess and refine the model-generated responses based on the given indoor or outdoor images and scenarios. You will evaluate the provided responses on several criteria and make necessary corrections to ensure accuracy, usability, and relevance. Each scenario consists of the following:

1. Image: an indoor or outdoor environmental scene provided as visual context.
2. Request: a BLV mobility-related request or task from BLV users.
3. Step-by-step description: the vision-language model’s response to the request.

**Notes for Refinement:** The deep context responses often follow this format:

1. Scene Description: An overview of the environment, highlighting key landmarks.
2. Distance and Direction to the Goal: Clear directional and distance information to the target.
3. Obstacles to Watch For: Specific obstacles the user should be aware of.
4. Step-by-Step Directions: Detailed instructions for completing the task.

When a response does not follow this format, refine it accordingly. Copy the original response and correct errors, remove unnecessary details, or add missing information.

**Evaluation Criteria:** For each response, you will rate the following aspects on a Likert scale (1 to 5) or a binary scale and provide corrections where necessary:

#### 1. *Direction Accuracy*

- Definition: Are the directions provided (e.g., angle, clock direction, left or right) in the response accurate and aligned with the image?

- Ratings: 1: no accurate info at all, 2: 3 inaccurate info, 3: 2 inaccurate info, 4: 1 inaccurate info, 5: entirely accurate

#### 2. *Depth Accuracy*

- Definition: Are the depth levels (e.g., x meters, x steps) in the response correct and aligned with the image?

- Ratings: 1: no accurate info at all, 2: 3 inaccurate info, 3: 2 inaccurate info, 4: 1 inaccurate info, 5: entirely accurate

#### 3. *Response Quality* (Step-by-Step Order, Relevance, Hallucination and Repeatedness)

- Definition: Is the response free from unnecessary repetition, illogical order, irrelevance to the scenario, and hallucinations?

- Ratings: 0: no (the response is illogically ordered, and includes irrelevant or hallucinated details, or has repetitions), 1: yes

#### 4. *Safety and Actionability*

- Definition: Are the actions in the response safe and actionable for BLV users?

- Ratings: 1: no safe or actionable actions, 2: no safe actions but actionable actions, 3: safe actions but include non-actionable actions (e.g., watch xxx, see xxx, observe xxx, drive xxx, etc.), 4: includes consensual safety actions (e.g., watch out for cars when crossing the street.), 5: entirely safe and actionable

#### 5. *Sufficiency*

- Definition: Does the response provide all the necessary information for BLV users to complete the task?

- Ratings: 1: strongly disagree, 2: disagree, 3: neutral, 4: agree, 5: strongly agree

#### 6. *Conciseness*

- Definition: Is the response concise and free from verbosity?

- Ratings: 1: strongly disagree, 2: disagree, 3: neutral, 4: agree, 5: strongly agree

#### 7. *Overall*

- Definition: Rate the response overall.

- Ratings: 1: not useful at all, 2: needs significant improvement, 3: needs some improvement, 4: valid for BLV users, 5: very useful for BLV users.

Table 12: **Instruction guideline for human experimental procedure.** We provide an overview and detailed steps for the annotation and refinement processes. The human annotators were also encouraged to ask clarifying questions to the authors whenever needed.

1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

Here are BLV user's request and generated response. Please rate and improve the generated response.

Request: Cross the street using the pedestrian crosswalk in the 3 o'clock direction.

1) Scene Description: You are currently walking on a sidewalk near the 3 o'clock direction, and there is a pedestrian crosswalk in front of you. The street is relatively empty but has some traffic.

2) Main Obstacles: Watch out for cars moving nearby, including vans parked on the opposite side of the street and vehicles approaching from behind you.

3). Step-by-Step Actions:

- Step 1: Wait until all cars have stopped at the red light before crossing the street to ensure safety.
- Step 2: Cross the street using the pedestrian crosswalk, making sure to look both ways for any oncoming traffic.
- Step 3: Once across the street, continue walking along the sidewalk towards your destination.

Are the directions (e.g., x o'clock, right or left) in the response accurate and aligned with the image?

no accurate info at all<sup>[1]</sup>  3 inaccurate info<sup>[2]</sup>  2 inaccurate info<sup>[3]</sup>  1 inaccurate info<sup>[4]</sup>

entirely accurate<sup>[5]</sup>

Are the depth levels (e.g., x meters, x steps) in the response accurate and aligned with the image?

no accurate info at all<sup>[6]</sup>  3 inaccurate info<sup>[7]</sup>  2 inaccurate info<sup>[8]</sup>  1 inaccurate info<sup>[9]</sup>

entirely accurate<sup>[0]</sup>

Is the response free from unnecessary repetition, illogical order, irrelevant to the scenario, and hallucinations?

yes<sup>[1]</sup>  no (the response is illogically ordered, includes irrelevant or hallucinated details, or has repetitions)<sup>[2]</sup>

Are the actions included in the response safe and actionable to BLV users?

no safe or actionable actions<sup>[3]</sup>  no safe actions but includes actionable actions<sup>[4]</sup>

safe actions but includes non-actionable actions<sup>[5]</sup>

includes consensual safety actions (e.g., watch out for cars when crossing the street)<sup>[6]</sup>

entirely safe and actionable<sup>[7]</sup>

Does the response provide all the necessary information for BLV users to fulfill their request?

Strongly Agree<sup>[8]</sup>  Agree<sup>[9]</sup>  Neutral<sup>[0]</sup>  Disagree<sup>[1]</sup>  Strongly Disagree<sup>[2]</sup>

Is the response concise and free from unnecessary verbosity?

Strongly Agree<sup>[3]</sup>  Agree<sup>[4]</sup>  Neutral<sup>[5]</sup>  Disagree<sup>[6]</sup>  Strongly Disagree<sup>[7]</sup>

Does the response provide all the necessary information for BLV users to complete the task?

not useful at all<sup>[8]</sup>  needs significant improvement<sup>[9]</sup>  needs some improvement<sup>[0]</sup>  useful for BLV users<sup>[1]</sup>

very useful for BLV users<sup>[2]</sup>

Please 'copy and paste' the response and 'revise' any wrong direction/depth information. 'Delete' any repetitions and irrelevant sentences. 'Reorder' the sentence if necessary. 'Add' more additional sentence if you answered the response is insufficient. Make sure to write with the proper format (1) Scene Description: xxx, ..., 3) Step-by-Step Actions: xxx

Enter your improved deep contexts here and press the Add button.

Add

Figure 4: Sample screenshot of interface used in the human experiment. This annotation screen with a different image-request-response is shown 100 times per annotator.

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

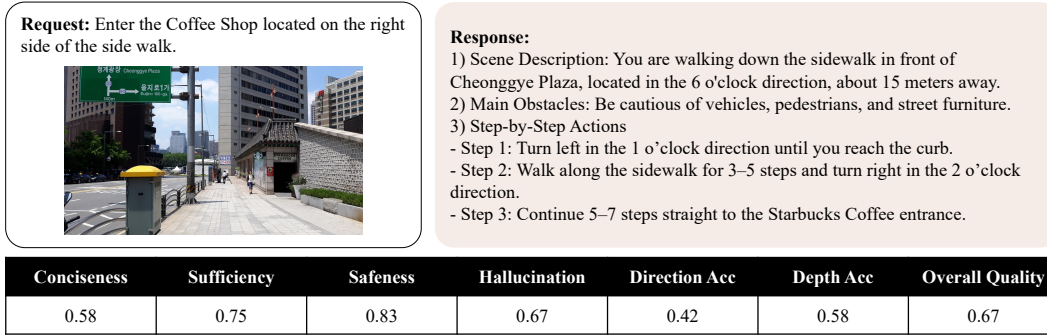


Figure 5: **Example of EYE4ALLMulti**. EYE4ALLMulti comprises a text request, an image, model-generated responses, and scores across seven dimensions: Conciseness, Sufficiency, Safeness, Hallucination, Direction Accuracy, Depth Accuracy, and Overall Quality. These scores are normalized and averaged over 2–3 human annotators.

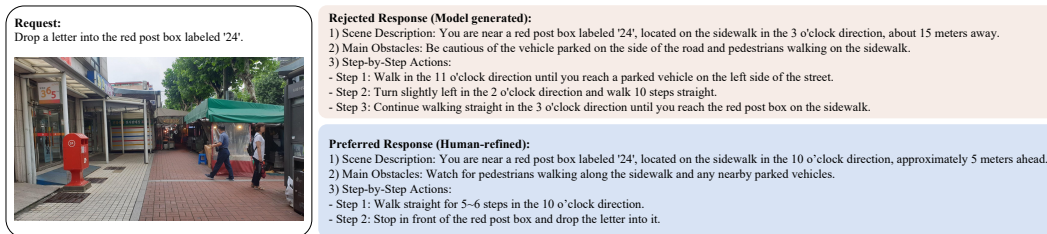


Figure 6: **Example of EYE4ALLPref**. EYE4ALLPref consists of a text request, an image, model-generated responses, and human-refined responses. For evaluation, model-generated responses are treated as rejected, while human-refined responses are considered preferred samples.

	PASCAL-50S P-Acc	FOILR1 P-Acc	FOILR4 P-Acc	FlickrExp $\tau_c$	FlickrCF $\tau_b$	Polaris $\tau_c$	Polaris* P-Acc	OID* P-Acc	ImgREW P-Acc
<b>MULTI-TAP</b>									
- Qwen-2B-S	84.7	97.3	97.3	57.8	39.4	59.8	83.3	58.5	59.9
- Qwen-7B-S	83.8	97.0	97.0	56.3	39.9	61.5	84.7	57.3	60.3
- InternLM-7B-S	82.0	97.1	97.1	53.1	38.8	57.3	83.7	58.9	54.3
- LLaMA-3.2-S	82.7	96.5	96.5	56.9	38.3	60.9	81.3	56.5	63.5

Table 13: **Performances of our MULTI-TAP trained with the empty prompt setting on 8 human preference judgment datasets**. The performances are similar to the main results in Table 1.

1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

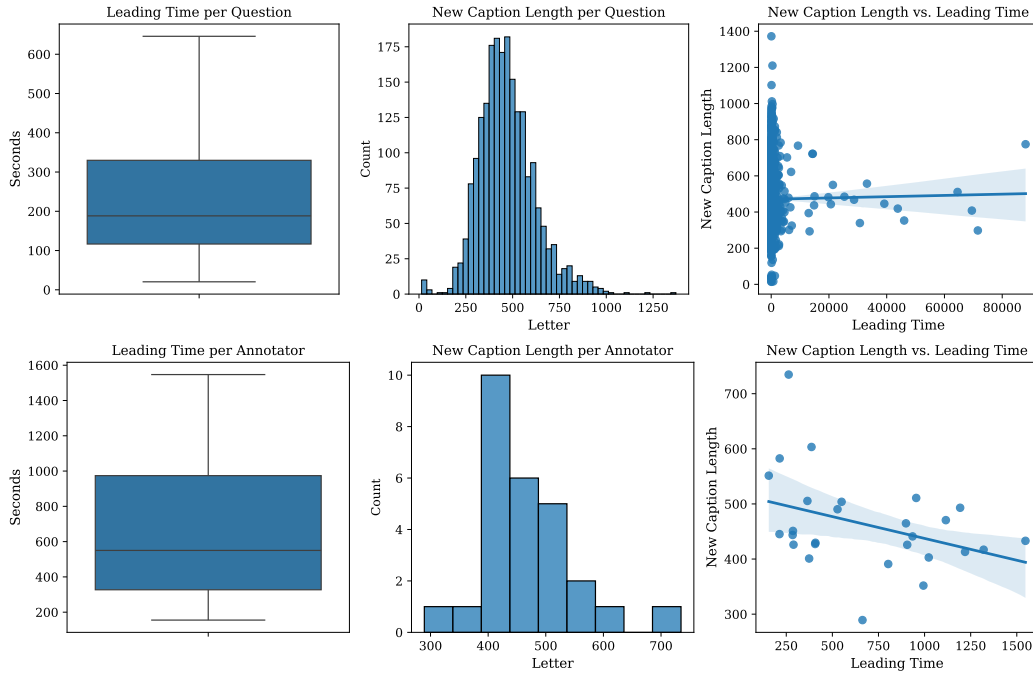


Figure 7: Distribution and correlation of leading time per question/annotator and the length of newly human-generated captions. Longer leading time does not necessarily mean more captions generated by human annotators.

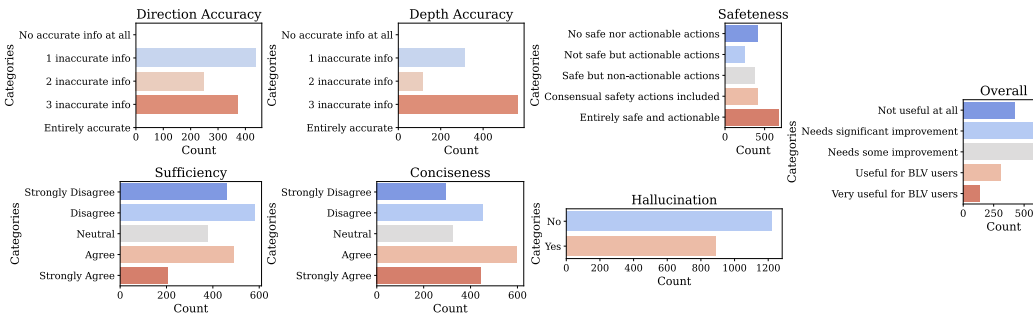


Figure 8: The distribution of human-annotated scores across seven categories. These plots indicate a room for significant improvement in LVLMs, especially in terms of accuracy, safety, and sufficiency.

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

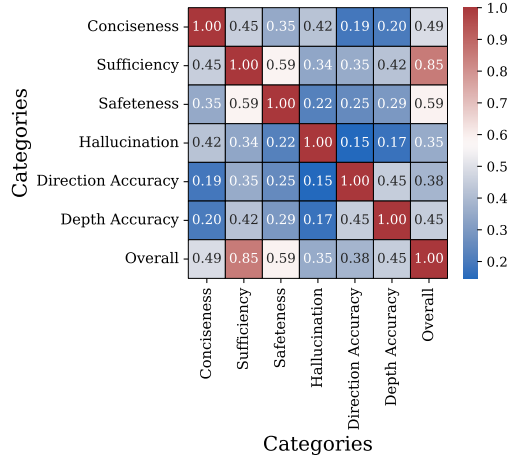


Figure 9: **Correlation of human judgment scores across different category pairs.** The overall quality of LVLm responses is highly correlated with the sufficiency category, captured with Pearson’s correlation coefficient of 0.85.

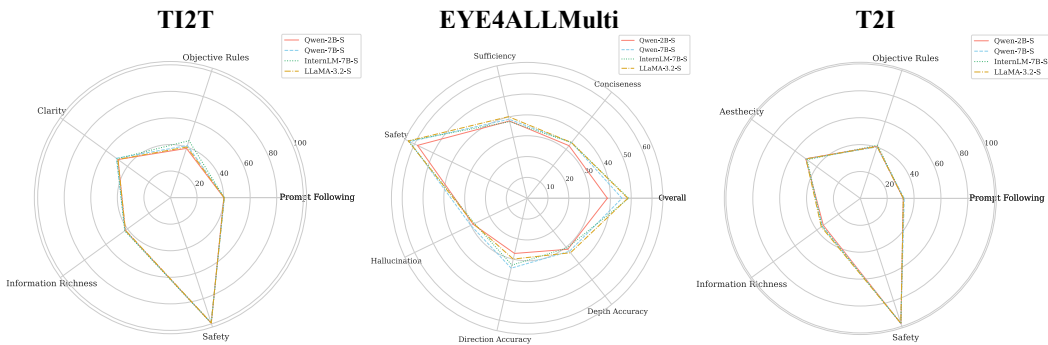


Figure 10: **Performances of our MULTI-TAP models on multi-objective scoring datasets.** The MULTI-TAP models trained on different LVLm architectures show a similar trend across multiple categories for all three datasets.

1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349

<i>Dimension</i>	VisionREW-S	MULTI-TAP <sub>Qwen-2B-S</sub>	MULTI-TAP <sub>Qwen-7B-S</sub>	MULTI-TAP <sub>InternLM-7B-S</sub>	MULTI-TAP <sub>LLaMA-3.2-S</sub>
0	90.38	91.85	90.87	91.76	90.78
1	1.08	99.41	99.12	99.41	98.92
2	94.70	97.25	96.47	97.25	96.17
3	10.50	94.80	94.01	94.80	93.92
4	26.20	75.76	75.17	73.21	74.88
5	10.79	93.72	92.35	93.33	92.84
6	99.90	99.80	99.51	99.80	99.71
7	93.62	93.62	92.44	93.52	92.54
8	11.09	90.19	89.50	90.19	89.89
9	1.77	99.41	98.92	99.41	99.12
10	99.12	99.51	99.21	99.51	99.21
11	92.54	95.49	94.80	95.49	94.31
12	14.23	89.21	88.13	88.42	87.44
13	98.14	96.86	96.57	96.86	96.37
14	78.90	73.01	69.09	70.26	70.17
15	7.75	89.99	89.89	89.79	89.40
16	0.39	99.41	99.61	99.41	99.61
17	93.52	95.68	94.41	95.39	94.60
18	11.48	91.56	90.97	91.36	90.38
19	92.15	95.98	95.29	95.39	94.70
20	5.50	94.11	93.82	94.01	93.42
21	98.72	99.61	99.31	99.61	98.92
22	74.39	81.55	79.39	79.29	79.29
23	9.22	85.67	84.69	85.38	84.89
24	99.61	99.80	99.61	99.80	99.41
25	89.01	92.74	91.66	92.64	91.76
26	99.90	100.00	100.00	100.00	100.00
27	94.80	92.93	92.44	92.93	92.15
28	7.46	92.35	91.36	92.25	91.46
29	0.29	99.21	99.21	99.21	99.02
30	99.41	99.12	98.72	99.12	98.63
31	89.50	90.19	89.30	89.89	88.52
32	34.54	62.90	61.63	61.04	62.90
33	7.75	90.87	90.28	90.97	90.28
34	2.75	97.94	97.35	97.94	97.84
35	0.20	100.00	99.80	100.00	99.80
36	94.50	93.72	93.03	93.72	92.84
37	41.41	55.94	56.33	54.47	53.29
38	8.93	90.68	89.79	90.78	89.79
39	0.98	98.72	98.72	98.72	98.33
40	97.35	93.33	92.74	93.33	92.54
41	69.58	68.99	66.34	65.55	65.55
42	54.47	53.97	54.96	50.83	52.99
43	46.91	53.97	53.09	51.62	53.58
44	45.44	55.45	52.89	51.72	54.37
45	99.71	99.41	99.21	99.41	99.21
46	96.57	94.80	94.21	94.80	94.21
47	81.06	73.80	72.13	72.72	71.74
48	59.27	55.62	55.64	55.94	56.33
49	54.96	53.97	53.48	53.88	56.43
50	99.90	99.71	99.61	99.71	99.51
51	96.66	96.57	96.07	96.57	96.07
52	89.30	87.83	86.85	87.63	87.14
53	73.21	71.34	69.87	68.79	68.89
54	70.36	68.20	66.44	66.54	65.85
55	13.35	93.03	91.85	92.93	91.95
56	7.16	94.21	93.42	94.21	93.62
57	3.93	96.57	96.07	96.67	96.07
58	0.20	99.80	99.80	99.80	99.71
<b>Avg</b>	53.30	88.00	87.30	87.40	87.20
<b>Time</b>	51 days	4 hrs	6 hrs	6 hrs	11 hrs

Table 14: **Multi-objective performances of VisionREW-S and our MULTI-TAP models on VisionREW dataset across 59 dimensions.** MULTI-TAP outperforms VisionREW-S in terms of both efficiency and performance.

1350

1351

1352

1353

<i>Dimension</i>	VisionREW-S	MULTI-TAP <sub>Qwen-2B-S</sub>	MULTI-TAP <sub>Qwen-7B-S</sub>	MULTI-TAP <sub>InternLM-7B-S</sub>	MULTI-TAP <sub>LLaMA-3.2-S</sub>
Conciseness	33.55	46.54	55.59	62.30	58.89
Sufficiency	54.21	43.45	46.33	46.96	47.50
Safety	26.20	72.95	73.70	73.80	73.80
Hallucination	53.14	42.81	45.26	46.65	47.92
Direction Acc	55.38	42.07	43.13	44.73	45.69
Depth Acc	51.01	41.00	47.07	48.35	44.62
Overall	59.96	39.62	42.07	41.75	42.71
<b>Avg</b>	47.63	46.92	50.45	52.08	51.59

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

<i>Dimension</i>	VisionREW-S	MULTI-TAP <sub>Qwen-2B-S</sub>	MULTI-TAP <sub>Qwen-7B-S</sub>	MULTI-TAP <sub>InternLM-7B-S</sub>	MULTI-TAP <sub>LLaMA-3.2-S</sub>
Prompt Following Rate	41.75	59.35	59.52	56.84	60.40
Objective Rules	15.80	84.20	84.11	84.23	83.47
Aestheticity	19.97	79.96	80.79	80.37	81.01
Information Richness	41.94	59.28	59.20	57.25	59.94
Safety	0.78	99.22	99.22	99.22	99.10
<b>Avg</b>	24.05	76.40	76.57	75.58	76.78

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

Table 16: **Multi-objective performances of VisionREW-S and our MULTI-TAP models on Align-anything (T2I)-Binary dataset across five dimensions.** MULTI-TAP significantly surpasses VisionREW-S performance on the T2I dataset.

<i>Dimension</i>	VisionREW-S	MULTI-TAP <sub>Qwen-2B-S</sub>	MULTI-TAP <sub>Qwen-7B-S</sub>	MULTI-TAP <sub>InternLM-7B-S</sub>	MULTI-TAP <sub>LLaMA-3.2-S</sub>
Prompt Following Rate	8.08	91.76	91.35	91.71	91.50
Objective Rules	4.72	94.82	94.77	95.18	95.08
Clarity	1.40	98.65	98.55	98.60	98.60
Information Richness	13.11	86.58	86.17	86.58	86.37
Safety	0.00	100.00	99.53	100.00	99.79
<b>Avg</b>	5.46	94.36	94.07	94.41	94.27

Table 17: **Multi-objective performances of VisionREW-S and our MULTI-TAP models on Align-anything (TI2T)-Binary dataset across five dimensions.** MULTI-TAP significantly surpasses VisionREW-S performance on the TI2T dataset.

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

Generated Samples Using Generative Model	
1. Accuracy:	Answer 1: Inaccurate. The image doesn't show a blond woman in a blue shirt waiting for a ride. Answer 2: More accurate. It correctly describes a young child in blue goggles in a pool setting.
2. Completeness:	Answer 1: Incomplete. It lacks key details about the child's position or the pool setting. Answer 2: More complete. It fully describes the child, their attire, and the surrounding environment.
3. Clarity:	Answer 1: Poor. It's vague and doesn't clearly convey what's in the image. Answer 2: Better. It's concise and easy to understand, directly describing what's visible.
4. Relevance:	Answer 1: Irrelevant. It doesn't address the question about the woman in the blue shirt. Answer 2: Relevant. It directly answers the question about the child in blue goggles.
Overall Judgment:	Answer 2 is better.
It accurately describes the key elements of the image, including the child, their attire, and the pool setting, which directly answers the question about what's in the picture.	
-----	
I'll evaluate the caption based on the given criteria:	
1. Accuracy:	The caption accurately describes the child wearing blue goggles and sitting in a float, which aligns with the visual information provided.
2. Completeness:	The caption fully addresses the question by mentioning both the child and the goggles/float, providing a comprehensive answer.
3. Clarity:	The caption is easy to understand, using simple language to convey the information clearly.
4. Relevance:	The caption directly relates to the question and the image, focusing on the key elements described in the visual content.
Overall Judgment:	0.9
The caption scores high in most categories, with perfect accuracy, completeness, and clarity. While it could be slightly more concise, it effectively captures the essential elements of the image. The only minor reason for a point deduction is that it could potentially be even more streamlined, but this doesn't significantly impact its overall quality.	

Table 18: **Samples of generated responses from Molmo-7B with the first (above) and second (below) prompts in Table 6.** Molmo-7B produces answers based on the criteria given in the prompt and provides the summarized reason for the overall judgment.

Generated Samples Using Generative Model	
1458	
1459	
1460	Evaluation:
1461	
1462	1. <b>Accuracy</b> :
1463	- <b>Answer 1</b> : This answer is accurate as it describes a city bus driving down a city street, which is clearly visible in the image.
1464	- <b>Answer 2</b> : While the sky is indeed blue and clear in the image, the description doesn't capture the primary aspects of the scene, such as the presence of the bus and other elements like the bank and parked cars.
1465	
1466	
1467	2. <b>Completeness</b> :
1468	- <b>Answer 1</b> : This answer covers more parts of the scene. It includes the bus, a critical element, and implicitly communicates the urban setting.
1469	- <b>Answer 2</b> : This answer focuses only on the sky and doesn't account for other significant scene elements.
1470	
1471	3. <b>Clarity</b> :
1472	- <b>Answer 1</b> : This answer is straightforward and easily understandable.
1473	- <b>Answer 2</b> : It's clear but less descriptive and not as comprehensive.
1474	
1475	4. <b>Relevance</b> :
1476	- <b>Answer 1</b> : Relevant as it directly mentions a city bus and the city street setting, aligning closely with the image.
1477	- <b>Answer 2</b> : While relevant to the image, the focus on the sky doesn't capture the main elements of the scene.
1478	
1479	Overall Judgment: <i>Answer 1 is better</i> .
1480	While both answers are relevant, Answer 1 captures the primary subjects of the image (city bus driving down a street) and provides a more comprehensive description. It is more accurate and addresses more aspects of the scene, making it the better choice for the given question.
1481	
1482	
1483	
1484	
1485	Evaluation:
1486	
1487	1. Accuracy:
1488	The caption generally aligns well with the visual elements of the image. It accurately captures the sinister and eerie feeling of the portal, the detailed look, and the dark tones. However, it doesn't mention the arch and pipework leading to the portal, nor does it touch on the light leak in the top part of the portal.
1489	
1490	2. Completeness:
1491	The caption provides a comprehensive description of the image, capturing many details like the ultra-detailed nature, the warm interior light, and the cinematic shot. However, it misses out on describing the arch and pipework that serve as pathways into the portal.
1492	
1493	
1494	
1495	3. Clarity:
1496	The caption is clear and uses descriptive language to convey the mood and tone of the image effectively. The terms like 'cinematic shot', 'ultra-detailed', and 'photo-realistic' are used to paint a vivid picture of the scene, making it easy for the reader to grasp the image's essence.
1497	
1498	
1499	4. Relevance:
1500	The caption is very relevant as it directly addresses the nature of the image and uses terminology that reflects the detailed and realistic rendering expected from tools like MCBess.
1501	
1502	Overall Judgment:
1503	Given these assessments, the caption is extremely good. It accurately reflects the contents of the image while using vivid, relevant, and descriptive language. However, it could be improved by noting additional observations, such as the arch and pipework leading into the portal, and the top light leak in the portal.
1504	
1505	
1506	Score:
1507	Overall Judgment: 0.85
1508	
1509	
1510	
1511	

Table 19: **Samples of generated responses from InternVL2-8B with the first (above) and second (below) prompts in Table 6.** InternVL2-8B answers based on the criteria given in the prompt and also provides the summarized reason for the overall judgment.

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

Generated Samples Using Generative Model	
The first caption, "A city bus driving down a city street," is better than the second caption, "the sky is clear and blue." The first caption accurately describes the main subject of the image, which is a city bus driving down a city street. This provides a clear and relevant description of the scene depicted in the image. On the other hand, the second caption, while true, does not capture the primary focus of the image and lacks context about the urban environment and the presence of the bus. Therefore, the first caption is more accurate and informative, making it the better choice.	
The better caption: [first].	
-----	
The LMM response provides a detailed description of the image, highlighting several key elements such as "deathly portal to the abyss," "ultra detailed," "warm interior light," "cinematic shot," "photorealistic," "octane render," "high definition," "fine details," "sinister tones," "8 k," and "mcbess mood."	
However, the response is somewhat disjointed and lacks coherence. It mentions multiple aspects without clearly connecting them to form a cohesive description of the image. For instance, while it describes the image as having a "deathly portal to the abyss," it also mentions "warm interior light" and "cinematic shot," which are contradictory in terms of atmosphere.	
Additionally, some terms like "8 k" and "mcbess mood" are not standard descriptive terms for images and may confuse the reader. The response could be improved by providing a more unified and clear description that captures the overall mood and elements of the image effectively.	
Score: 60	
Reasons:	
- The response includes several relevant details about the image.	
- However, it is somewhat disorganized and lacks coherence.	
- Some terms are not standard or clear, which might confuse the reader.	
- Overall, the response could be more polished and focused.	

Table 20: **Samples of generated responses from LLaVA-Critic-7B with the first (above) and second (below) prompts in Table 6.** LLaVA-Critic-7B presents multiple justifications for its selection of the preferred caption and assigns an overall score accordingly.

	PASCAL-50S P-Acc	FlickrExp $\tau_c$	FlickrCF $\tau_b$	Polaris* P-Acc	OID* P-Acc	ImgREW P-Acc
IXCREW-S-v1	73.9	21.6	25.7	<b>89.1</b>	56.9	53.6
IXCREW-S-v2	<b>76.2</b>	<b>25.6</b>	<b>28.0</b>	88.7	<b>58.5</b>	<b>56.3</b>
IXCREW-S-v3	74.2	17.0	25.7	81.9	57.5	53.6

Table 21: **Performances of IXCREW-S using two different prompts in Table 9 and the same prompt as ours (with no content for the user).** Although version 3 is reported in the main paper to match the prompt setting we use for our models, here, we show that slightly different performances can be achieved with different prompt settings.

	FlickrExp $\tau_c$	FlickrCF $\tau_b$	ImgREW P-Acc
LLaMA-3.2-11B-v1	-7.88	5.49	46.0
LLaMA-3.2-11B-v2	<b>5.29</b>	<b>9.00</b>	<b>51.6</b>

Table 22: **Performances LLaMA-3.2-11B using two different prompts in Table 8.** In the main paper, we report version 2, which shows higher generalizability across datasets.