

Supplementary Materials of Instruction-Oriented Preference Alignment for Enhancing Multi-Modal Comprehension Capability of MLLMs

1. More Implementation Details

1.1. Training Configuration

For model optimization, we leverage Direct Preference Optimization (DPO) with Low-Rank Adaptation (LoRA). Specifically, we use the AdamW optimizer with a cosine learning rate scheduler, starting from the initial learning rate of $5e-6$. The total batch size is set to 8 and the gradient accumulation step is 8, resulting in an effective batch size of 64. All the models are trained for 1 epoch in the comparisons.

1.2. Prompt Design

We detail the prompts used in the automated preference construction. Specifically, these comprises the prompt I^{fb} for feedback generation (Equation 3), the prompt I^{rev} for response refinement (Equation 4), and the prompt I^{ver} for instruction-oriented verification (Equation 5). The prompts are displayed in Table 2, 3 and 4 respectively.

2. Qualitative Analysis

As illustrated in Figure 1, we analyze how each stage of our automated preference construction contributes to the construction of high-quality preference pairs. The revision stage improves the sampled responses by correcting factual errors and expanding the depth of the responses. Through critical feedback generation and refinement, potential flaws have the opportunity to be corrected and the responses are enriched with more valuable information. The verification stage then quantifies instruction fulfillment capacity of the responses. Responses which pass the verification are iden-

tified as preferred due to their stronger intrinsic alignment with instruction objectives.

As demonstrated in Figure 2, we present comparative examples to illustrate the improvements achieved by our approach over the baseline model. The visual comparisons reveal that IPA-aligned responses exhibit enhanced capability in addressing the intrinsic demands of instructions. These examples validate that our approach enables the model to better interpret and fulfill the requirements embedded in diverse instructions, moving toward more comprehensive understanding capability.

3. More Experiments

3.1. Experiments with More Advanced Models

We also conduct experiments using larger and more advanced models, InternVL3-14B and Qwen2.5VL-72B, where the same dataset as in the main manuscript is used for training these models. As shown in Table 1, despite the preference dataset being constructed with a relatively weaker and smaller model (Qwen2VL-7B), it brings consistent gains to significantly larger models, demonstrating strong generalizability across model scales and architectures. While the improvements are more moderate due to the already strong baselines, our method still yields meaningful enhancements in overall comprehension ability. We believe that using stronger models for data collection can yield further improvements.

Table 1. Experiments with more advanced models. The training dataset is the same as in the main manuscript.

| Model | Hallucination | | | General VQA | | | | | Text | |
|----------------|--------------------------|-----------|----------|---------------------|-----------|------------------------|------------------|------------|----------|-----------|
| | HallBench _{avg} | POPE | MMHal | MMMU _{val} | MMStar | MMVet _{turbo} | MME _P | LLaVABench | OCRBench | |
| InternVL3-14B† | 55.7 | 89.5 | 3.8 | 63.1 | 68.8 | 73.0 | 1740.5 | 83.0 | | 87.6 |
| + IPA (ours) | 56.0 +0.3 | 89.6 +0.1 | 3.9 +0.1 | 63.3 +0.2 | 69.1 +0.3 | 75.3 +2.3 | 1742.3 +1.8 | 89.9 +6.9 | | 88.0+0.4 |
| Qwen2.5VL-72B† | 55.3 | 86.4 | 4.3 | 68.1 | 71.6 | 74.4 | 1721.3 | 97.6 | | 88.8 |
| + IPA (ours) | 58.4 +3.1 | 86.9 +0.5 | 4.4 +0.1 | 69.0 +0.9 | 71.5 -0.1 | 74.6 +0.2 | 1730.8 +9.5 | 102.2 +4.6 | | 89.2 +0.4 |

Table 2. Prompt I^{fb} for feedback generation.

You are a evaluation model tasked with providing a detailed critical analysis of the given instruction-response pair. Please evaluate the response based on the criteria below:

Accuracy of Image Content Description

- Assess whether all descriptions of the image content in the response are accurate.
- Identify any instances of hallucinations, where the response describes elements inconsistent with the image.
- If any inaccuracies or hallucinations are found, provide a corrected description.

Correctness of Knowledge

- Verify the factual accuracy of any commonsense statements or knowledge presented in the response.
- If any incorrect knowledge or misconceptions are identified, provide the correct information.

Validity of Reasoning and Inferences

- Evaluate the logical consistency of the reasoning processes and inferences made in the response.
- Ensure that conclusions are appropriately derived from correct premises.
- If any flawed reasoning or incorrect inferences are detected, provide a revised reasoning or inference.

Verifiability and Expression of Uncertainty

- Determine if the response includes statements that are difficult to verify or falsify.
- Check whether the response clearly indicates uncertainty for such statements when appropriate.
- If unverifiable statements are present without indicated uncertainty, suggest how to express the uncertainty clearly.

Adherence to Instructions

- Note any deviations from the instructions or omissions of important requirements.
- If the response deviates from the instructions or omits key elements, provide a corrected version that fully adheres to the instructions.

You should offer feedback detailing all identified issues, followed by a comprehensive evaluation and actionable steps for improvement with direct corrections.

Output Format:

Accuracy of Image Content Description

...

Correctness of Knowledge

...

Validity of Reasoning and Inferences

...

Verifiability and Expression of Uncertainty

...

Adherence to Instructions

...

Overall Assessment

...

Guidance for Improvement

...

Input:

<Instruction>{instruction}</Instruction>

<Response>{response}</Response>

Table 3. Prompt I^{rev} for response refinement.

Given an instruction-response pair related to an image, an assistant has provided feedback identifying issues in the response. This feedback enables the generation of an improved response.

So your task is to derive a revised response by:
Response + Feedback \rightarrow Revised Response

Output Format:

Revised Response

...

Input:

<Instruction>{instruction}</Instruction>

<Response>{response}</Response>

<Feedback>{feedback}</Feedback>

Table 4. Prompt I^{ver} for instruction-oriented verification.

You are an evaluation model. Given an instruction-response pair related to an image, your task is to make evaluation by assessing the alignment between the response and a provided reference response.

You can assess the alignment by following the criteria below:

- Can the reference answer be inferred from the response, i.e., Response \implies Reference?
- Are there critical details from the reference missing in the response?
- Are there contradictions between the response and the reference?
- Does the response adhere to the instruction?

Note:

- It is permissible for the response to contain additional information if it is related to the instruction and not contradictory to the reference.
- If this is a multiple-choice/yes-no question, does the response clearly provide the answer?

You should indicate whether the response pass the evaluation (“Yes” or “No”) concerning all the criteria.

Output Format:

Assessment

...

Verification

(Yes or No)

Input:

<Instruction>{instruction}</Instruction>

<Response>{response}</Response>

<Reference>{reference response}</Reference>

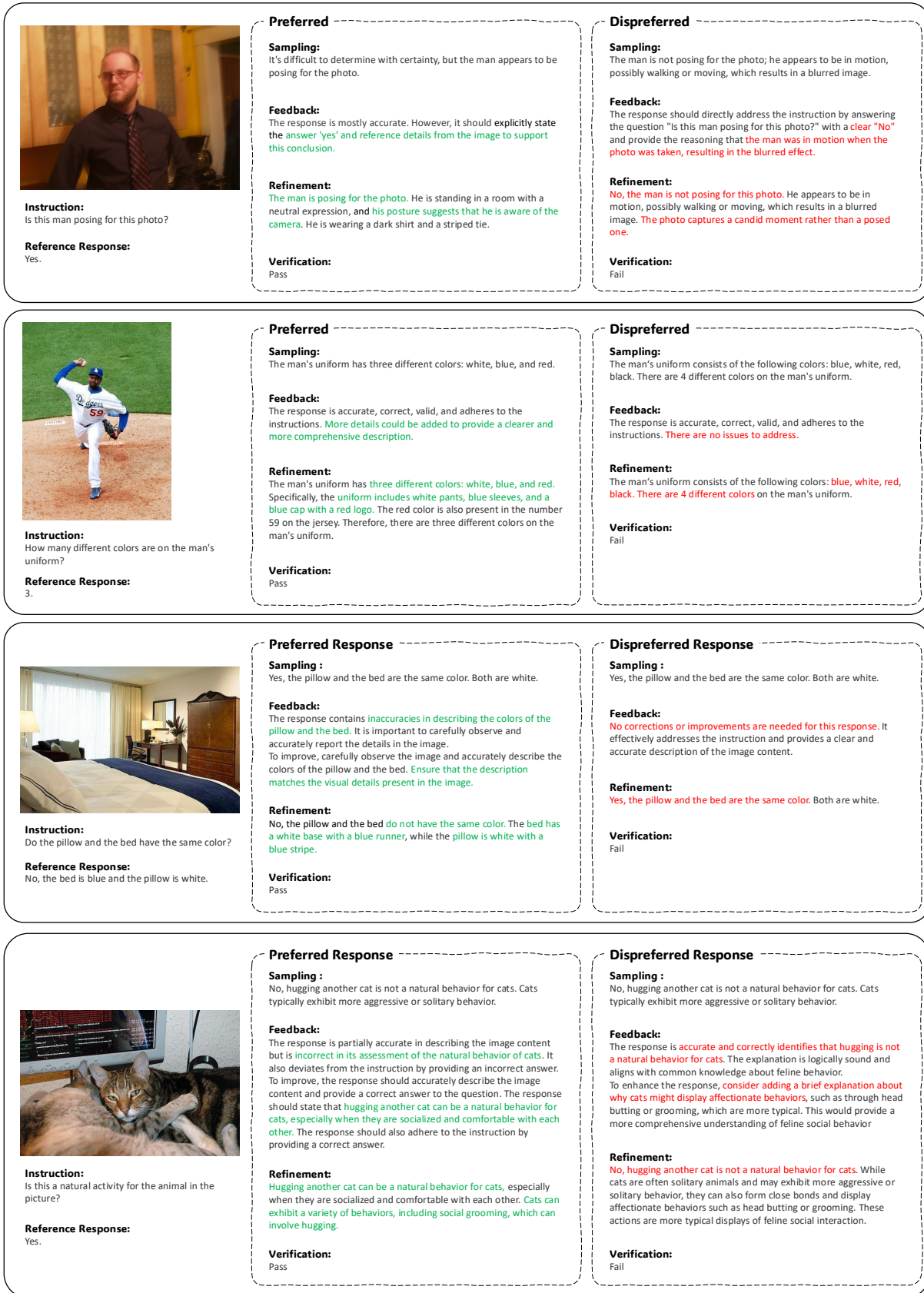


Figure 1. Some examples of the automated preference construction process.

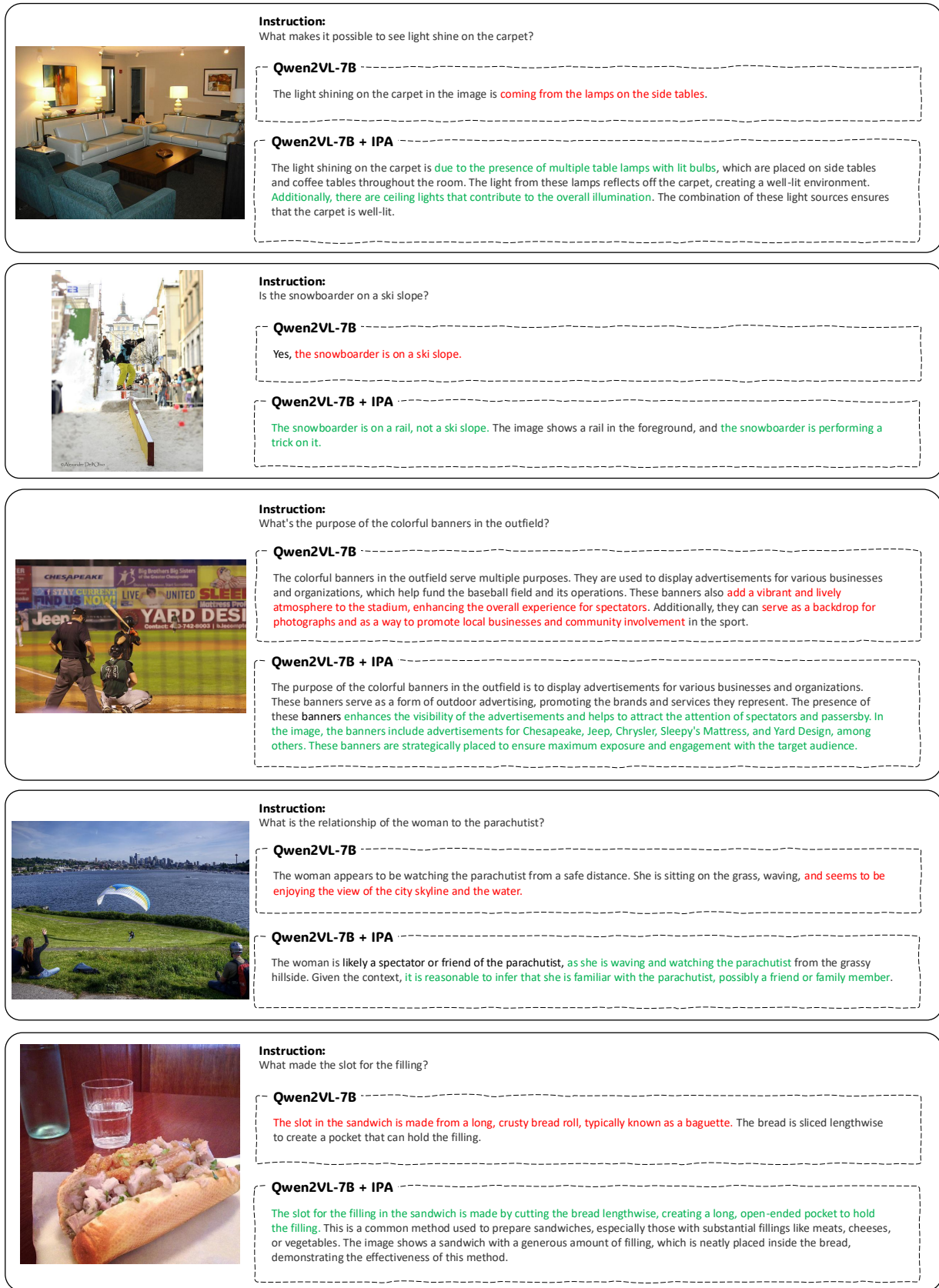


Figure 2. Some examples of the comparison between the baseline model the the model aligned with our approach.