# Towards Open-ended Visual Quality Comparison

Haoning Wu[⋆1], Hanwei Zhu[⋆2], Zicheng Zhang[⋆3], Erli Zhang[1]
Chaofeng Chen[1], Liang Liao[1], Chunyi Li[3], Annan Wang[1], Wenxiu Sun[4]
Qiong Yan[4], Xiaohong Liu[3], Guangtao Zhai[3], Shiqi Wang[2], and Weisi Lin[1]

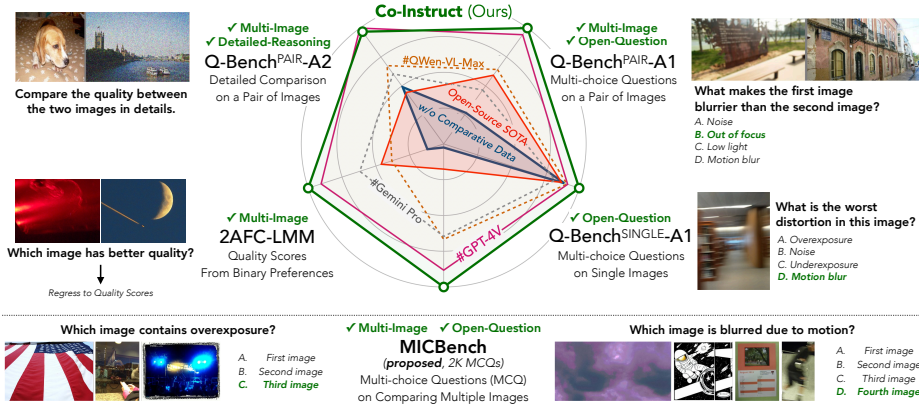https://huggingface.co/q-future/co-instruct



**Fig. 1:** The proposed **Co-Instruct**, *first-of-its-kind* open-source LMM with capability on *open-question* & *detailed-reasoning* visual quality comparison. It outperforms existing LMMs on the proposed **MICBench** as well as existing quality evaluation benchmarks.

**Abstract.** Comparative settings (*e.g. pairwise choice, listwise ranking*) have been adopted by a wide range of subjective studies for image quality assessment (IQA), as it inherently standardizes the evaluation criteria across different observers and offer more clear-cut responses. In this work, we extend the edge of emerging large multi-modality models (LMMs) to further advance visual quality comparison into open-ended settings, that **1)** can respond to *open-range questions* on quality comparison; **2)** can provide *detailed reasonings* beyond direct answers. To this end, we propose the **Co-Instruct**. To train this *first-of-its-kind* open-source open-ended visual quality comparer, we collect the Co-Instruct-562K dataset, from two sources: **(a)** LLM-merged single image quality description, **(b)** GPT-4V *"teacher"* responses on unlabeled data. Furthermore, to better evaluate this setting, we propose the **MICBench**, the first benchmark on multi-image comparison for LMMs. We demonstrate that **Co-Instruct** not only achieves in average **30%** higher accuracy than state-of-the-art open-source LMMs, but also outperforms GPT-4V (*its teacher*), on both existing related benchmarks and the proposed **MICBench**.

**Keywords:** Large Multi-modality Models (LMM) · Visual Quality Assessment · Visual Quality Comparison · Visual Question Answering
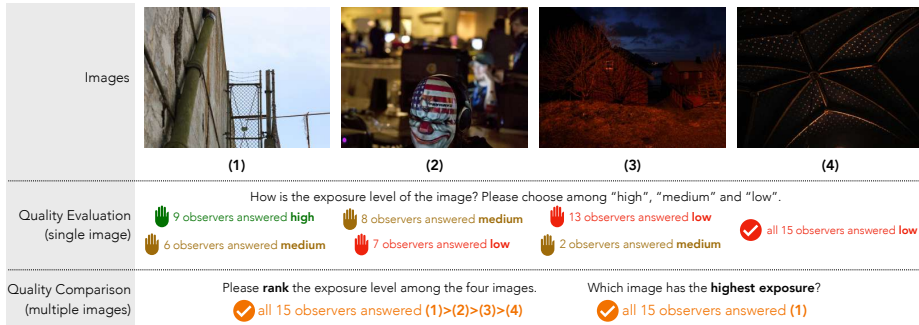
**Fig. 2: The motivation of** open-ended visual quality comparison: comparative settings can effectively avoid the **ambiguity on absolute evaluations** for single images, and provide more clear-cut judgements to serve as downstream guidances [53,65].

# 1  Introduction

Image quality assessment (IQA) has been an important domain in visual computing, as it provides effective recommendation [57] on high-quality visual contents and valuable guidance [53,61,65] for potential improvements. Most recently, several pioneer studies [49–51] have explored large multi-modality models (LMMs, *e.g.* GPT-4V) [4, 5, 9, 25, 34, 56], on expanding IQA from giving a scalar score (*e.g.* 3.457) to the open-ended scenarios, that allows evaluations in response to *open-range questions*, and provide *detailed reasonings* beyond an overall score.

While these studies sufficiently *emulate human ability* on IQA, they also suffer from the same drawback as human: **ambiguity on absolute evaluations**. For instance, as shown in Fig. 2 (a), different human observers hold different standards on the exposure level on single images, and henceforth provide diversified absolute evaluations. Nevertheless, while asked to compare the exposure level of the images, all observers agree with the rank (1)>(2)>(3)>(4) (Fig. 2(b)); all observers also agree that (1) has the highest exposure, though not all choose the option *high* while evaluating it independently. Given this observation, comparison has been a traditional human study setting for quality evaluation [43] and adopted by a wide range of existing subjective studies [10,35,52,53,62]. Furthermore, to avoid the ambiguity, the comparative settings are also predominantly adopted [53,65] while applying IQA algorithms for improvement guidance.

While comparison has widely-recognized significance for IQA, existing related datasets [35] and methods [28,63] are generally based on overall quality comparison and have not extended to the open-ended scenarios; on the other hand, open-source LMMs [4,25,50] are usually only fine-tuned with *single image* instruction tuning datasets [1,13,26] and proved to lack enough capability even on two image comparison settings [67,69]. While these gaps have clearly indicated the need of a specific instruction tuning dataset for visual quality comparison, it is too expensive to collect such a dataset from human. To avoid costly human labors, we propose an alternative strategy to collect the training dataset, named _Collaborative Instruction Tuning from Weak Supervisors_ (*Co-Instruct*). Specifically, we adopt
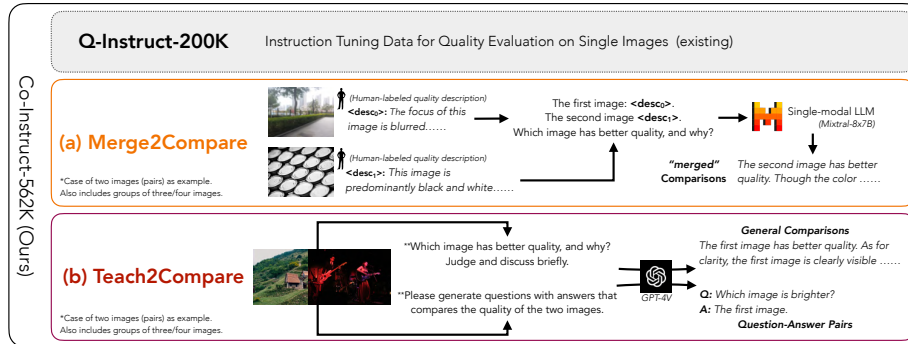
**Fig. 3:** **The construction methodology of** Co-Instruct-562K, a combination of **(a) Merge2Compare** (*LLM comparison from human-labeled single image quality descriptions*) and **(b) Teach2Compare** (*GPT-4V comparison on multiple unlabeled images*).

two non-perfect supervisors: **1) Merge2Compare** (Fig. 3(a)). Originated from single image quality descriptions on 19K images as labeled by human in the Q-Pathway [50] dataset, we randomly match them into 100K groups (*2-4 images per group*) with removing the most similar descriptions with an text embedding model [44]. Then, we prompt a single-modal large language model (LLM) [16] to compare multiple human descriptions in a group, and "*merge*" them into 100K pseudo comparisons. **2) Teach2Compare** (Fig. 3(b)). Observing that GPT-4V has especially high accuracy on pairwise settings [67,69] among existing LMMs, following existing practices [2,17], we leverage GPT-4V responses to further expand our dataset. We collect 9K unlabeled images and randomly match into 30K image groups (*also 2-4 images per group*), and obtain GPT-4V responses on both *caption-like* general comparisons and *question-answer pairs* for comparisons. By integrating Q-Instruct-200K [50] (*on single images*), **Merge2Compare**, and **Teach2Compare**, we construct the Co-Instruct-562K, the first instruction tuning dataset for open-ended visual quality comparison.

To correctly refer to each specific image during conversation, we define a specific image-text interleaved format [5] to handle multi-image cases, as follows:
`User: The first image: <img`$_0$`> The second image: <img`$_1$`> ... <query>`
`Assistant: <response>`

Moreover, as we need to feed multiple images together to the LLM decoder, adopting the most popular LLaVA [25, 26, 39] structure that linearly projects visual embeddings will exceed the context window of the language models [42] and cause errors. Henceforth, we adopt an alternative visual abstractor structure [56] to first reduce visual token length (*from 1,025 to 65 tokens per image*), and then concatenate them with text embeddings to pass to language decoders. By learning from the Co-Instruct-562K dataset and the specially-designed input structure, we present the **Co-Instruct**, with up to **86%** improvements than its baseline [56], and **61%** better than the state-of-the-art open-source LMM. More importantly, despite using GPT-4V as one of its teachers, it still surpasses the GPT-4V *teacher* in a variety of multi-choice question (MCQ) benchmarks, and matches GPT-4V ability in scenarios requiring detailed language reasonings.

After training the model **Co-Instruct**, our another concern is the lack of abundant evaluation settings on multi-image comparison: while Q-Bench [49,67] series have covered multiple formats on single images and image pairs, there is no existing evaluation scenario for quality comparison **beyond two images**. To complement this evaluation setting, we construct the **MICBench**. Specifically, the **MICBench** contains 2,000 multi-choice questions (MCQ) comparing the quality or related attributes among a group of three or four images (*each half*), in which over half of the questions are *Which* questions (*e.g. which image has highest clarity?*). Aiming to extract an image with a specific quality-related appearance from a group, *Which* questions are the most important questions related to image comparison. Despite *Which* questions, the **MICBench** also includes *Yes-or-No* questions and other question types (*What/How/How-Many, etc*) to provide a holistic benchmark setting for multi-image quality comparison.

In summary, we conduct a systematical study towards *open-ended* visual quality comparison. Our contributions can be summarized as three-fold:

1. We construct the first instruction-tuning **dataset** for visual quality comparison, the Co-Instruct-562K. With data collected from two "*weak supervisors*", **Merge2Compare** (LLM-merged comparisons) and **Teach2Compare** (GPT-4V pseudo-labeled comparisons), our public dataset significantly expands the capabilities of open-source LMMs on visual comparative settings.
2. We propose the most capable **model** for open-ended visual comparison, the **Co-Instruct**. With image-text interleaved input format and fine-tuned with the Co-Instruct-562K dataset, it significantly outperforms existing methods (and even GPT-4V) in multiple open-ended visual comparison tasks. With open weights, it allows for broader application than proprietary models.
3. We construct the **benchmark**, the **MICBench**, as the first benchmark to evaluate LMMs for quality comparison on multiple images (more than two). It covers 2,000 diverse-type open-range MCQs related to visual quality comparison among three or four images. The **MICBench** contributes to more holistic evaluation studies on the visual quality comparison problem.

## 2   Related Works

### 2.1   Visual Quality Comparison

Visual quality comparison (*especially* paired comparison) is a widely used subjective quality assessment methodology, serving as the most reliable way to collect human opinions [31]. However, when the number of images increases, the experiments become infeasible because of the exponential growth of pairwise comparisons [14]. While many active sampling methods have been proposed to reduce the number of pairs [21, 33, 54], they are computationally expensive and unsuitable for large-scale experiments. Despite subjective studies, learning to rank quality is widely proven as effective by many objective approaches [8,27,30,46,48,63,64]. Nonetheless, they typically only predict a scalar score or a binary judgement for overall comparison, limiting their ability to provide meaningful feedbacks into specific types of distortions or preferences.

## 2.2   LMMs for Visual Quality Evaluation

Several recent studies have explored LMMs for visual quality evaluation. The Q-Bench [49] proposes a holistic benchmark for LMMs on low-level perception (*quality-related MCQs*), description (*quality-related captioning*) and assessment (*predicting scores*). Following this path, Q-Instruct [50] advances the ability of LMMs with a large-scale human-annotated dataset, and Q-Align [51] designs a text-guided instruction tuning for score predictions and outperforms non-LMM approaches. However, these explorations are based on single images and have not covered comparative settings. While most recent benchmarks [67,69] suggest that open-source LMMs trained with single image datasets cannot perform well on comparative settings, to bridge this gap, we collect the first instruction tuning dataset to teach LMMs to compare visual quality, the Co-Instruct-562K, and our model significantly improves the ability of open-source LMMs on comparative settings, moving a step forward on the basis of existing explorations.

**Table 1:** Statistics of our Co-Instruct-562K dataset. '#' denotes *"the number of"*.

| Subsets | Q-Instruct-200K [50] | Merge2Compare | Teach2Compare-*general* | Teach2Compare-*Q&A* | All |
|---|---|---|---|---|---|
| Instruction Type | Detailed Reasoning, Question Answering | Detail Reasoning | Detail Reasoning | Question Answering | – |
| # Total Images | 19K *(shared, both using Q-Pathway images)* | | 9K *(using shared images)* | | 28K |
| # Total Data Items | 202K | 100K | 30K | 230K | 562K |
| # Single Images | 202K | 0 | 0 | 0 | 202K |
| # Image Pairs | 0 | 70K | 18K | 134K | 222K |
| # Groups of Three | 0 | 20K | 6K | 51K | 77K |
| # Groups of Four | 0 | 10K | 6K | 45K | 61K |

# 3   Data Construction

In this section, we elaborate on the construction process of the Co-Instruct-562K dataset. Though human annotation is the most direct approach to collect data, as is widely acknowledged [21,54], acquiring sufficient comparative data on a large set of images demands a markedly increased volume of human annotations [10, 35] in contrast to gathering opinions on the same set of individual images. To avoid this unbearable cost, we propose an alternative data construction strategy *without additional human annotation*, by following three key principles:

1. *Convert*: Utilize reliable information from existing datasets.
2. *Learn-from-Model*: Leverage the verified capabilities of models.
3. *Co-Instruct*: Collecting diverse subsets that complement each other.

Under the principles, we collect two different subsets for instruction tuning: **Merge2Compare** (Sec. 3.1), which *converts* information from human quality descriptions on single images, and the ability of single-modal LLMs on comparing and analyzing texts; and **Teach2Compare** (Sec. 3.2), which leverages the verified ability of GPT-4V on comparing images. Finally, we discuss how the two subsets complement each other (Sec. 3.3) under a *co-instruct* scheme.

## 3.1   Merge2Compare.

In this part, we define the construction process for **Merge2Compare** includes three steps: 1) pair/group matching; 2) top-similarity pair removal; and 3) LLM *merging*. An examplar illustration of the process is shown in Fig. 4.
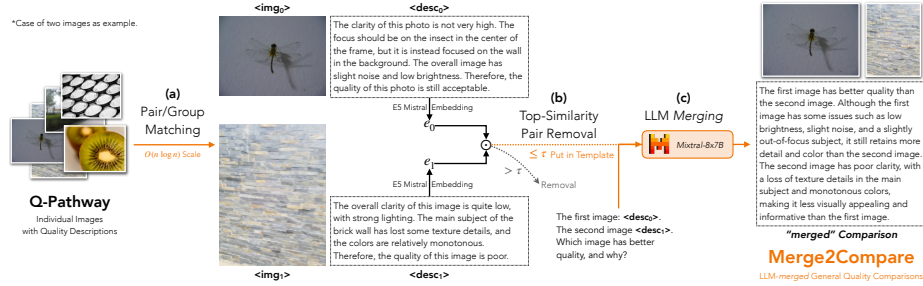
**Fig. 4:** The pipeline of constructing **Merge2Compare**: images are first matched into groups **(a)**, and then filtered via top-similarity removal **(b)**. After filtering, the single image quality descriptions are *merged* **(c)** into comparisons by the LLM [16].

*Step 1: Pair/Group Matching (Fig. 4(a)).* To best utilize information from existing single image quality descriptions, following the empirical rule to sample $O(n \log n)$ pairwise combinations to effectively rank among all individual items in a set [37, 45], we randomly sample 81K image pairs from all 19K images in the Q-Pathway. Despite pairs, we further sample 27K groups with three images and 18K groups with four images to cover the scenarios of more images.

*Step 2: Top-Similarity Pair Removal (Fig. 4(b)).* The effectiveness of the *merge* comes from the differences among descriptions, *e.g.* between *The quality is acceptable* for <img_0> and *The quality is poor* for <img_1>. However, if descriptions in a pair/group contains almost the same information (*e.g.* both images with *The clarity is good, but lighting is dark*), the *merged* comparisons will lack enough information or even with false predictions. Henceforth, we use E5-Mistral [44] text embedding model to compute similarities among descriptions, and remove if *any* high-similarity description pairs exist in the group. After removal, 70K image pairs (86% of initial samples), 20K groups of three (74% of initial) and 10K groups of four (55% of initial) are preserved and fed into the LLM for *merging*.

*Step 3: LLM Merging (Fig. 4(c)).* The key step for the **Merge2Compare** is to prompt LLMs to convert the single image evaluations to comparative texts. Specifically, following many existing practices [26,58,68], we put the descriptions as alternates of images in the context. Denote the description for image <img_i> as <desc_i>, the user query for LLM *merging* is formulated as follows:
(Pairs) `The first image: <desc`$_0$`> The second image: <desc`$_1$`>`
`Which image has better quality, and why?`
(Groups of Three/Four) $\{$`The `$K_{i+1}$` image: <desc`$_i$`> `$\mid_{i=0}^{N-1}\}$
`Please rank the quality of the images and justify your rankings.`
where $K_i$ represents the ordinal form of $i + 1$, *e.g.* $K_1$ is `first`, $K_2$ is `second`.

The *merged* comparisons are **overall comparisons with reasonings** (Fig. 4 *right*). To validate their reliability, we conducted a human examination on 250 random samples from **Merge2Compare** mixed with merged comparisons from 250 random removed groups. The correctness rate in **Merge2Compare** is 96%, while it is only 72% in the removed groups, demonstrating the effects of the top-similarity removal process in ensuring the quality of training data.
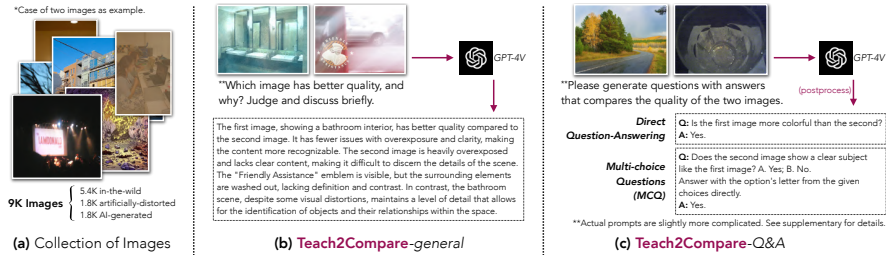
*Case of two images as example.

**9K Images**
- 5.4K in-the-wild
- 1.8K artificially-distorted
- 1.8K AI-generated

**(a) Collection of Images**

**"Which image has better quality, and why? Judge and discuss briefly.**

The first image, showing a bathroom interior, has better quality compared to the second image. It has fewer issues with overexposure and clarity, making the content more recognizable. The second image is heavily overexposed and lacks clear content, making it difficult to discern the details of the scene. The "Friendly Assistance" emblem is visible, but the surrounding elements are washed out, lacking definition and contrast. In contrast, the bathroom scene, despite some visual distortions, maintains a level of detail that allows for the identification of objects and their relationships within the space.

**(b) Teach2Compare**-*general*

**"Please generate questions with answers that compares the quality of the two images.**

*Direct Question-Answering*
**Q:** Is the first image more colorful than the second?
**A:** Yes.

*Multi-choice Questions (MCQ)*
**Q:** Does the second image show a clear subject like the first image? A. Yes, B. No.
Answer with the option's letter from the given choices directly.
**A:** Yes.

**Actual prompts are slightly more complicated. See supplementary for details.

**(c) Teach2Compare**-*Q&A*

**Fig. 5:** The pipeline of constructing **Teach2Compare**: 9K diverse images are collected and matched into 30K groups **(a)**. The groups are then fed to GPT-4V to obtain *general* quality comparisons **(b)** and *question-answering* **(c)** related to quality comparisons.

### 3.2   Teach2Compare.

Given existing evaluations [67, 69] suggesting that GPT-4V is decent at comparing visual quality (Tab. 2/3/4), we propose to collect GPT-4V responses as pseudo labels for quality comparison. As shown in Fig. 5, we collect diverse unlabeled images and feed them to GPT-4V with different prompts to obtain **Teach2Compare**-*general* (overall quality comparison) and **Teach2Compare**-*Q&A* (question-answer pairs related to quality comparison). Details as follows.

*Collection of Images (Fig. 5(a)).* For **Teach2Compare**, we collect 9K images from various sources to cover different quality concerns and visual appearances: **1)** 5.4K *in-the-wild* images from YFCC-100M [41] database; **2)** 1.8K images *artificially-distorted* images from COCO [3] (with 15 types of distortions via ImageCorruptions [32]) and KADIS-700K [24] (25 types of distortions); **3)** 1.8K *AI-generated* images from ImageRewardDB [53]. These 9K diverse unlabeled images are further grouped into 18K pairs, 6K groups of three, and 6K groups of four, for GPT-4V to provide pseudo labels under two formats, as follows.

**Teach2Compare**-*general (Fig. 5(b)).* Similar to **Merge2Compare**, the *general* subset also consists of overall comparison with reasonings. Specifically, we substitute the `<desc_i>` in the **Merge2Compare** prompt template to respective real images `<img_i>` to feed to GPT-4V. After collection, we also conduct a similar 250-sample spot check on the output pseudo labels, which reports around 94% correctness. Though slightly less accurate than **Merge2Compare** (96%), examiners also report that GPT-4V labels contain more **content information** which has been observed to enhance quality understandings of models [19, 20, 47]. The two subsets are expected to complement each other for better learning outcomes.

**Teach2Compare**-*Q&A (Fig. 5(c)).* Despite general comparisons, for GPT-4V, we also collect a specific subset to improve LMM ability on responding to open-range questions. To achieve this, we first list reference aspects (*clarity, lighting, color, etc*) and then ask GPT-4V to generate questions (and respective correct answers, and false answers for the questions) on comparing these aspects among the images. After removing failed generations, we obtain 230K question-answers from 30K image groups. These question-answers are converted to both direct question-answering and multi-choice questions (as A-OKVQA [38]) for training.
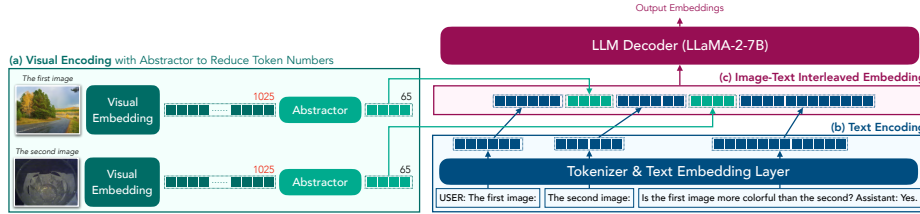
**Fig. 6: The structure of Co-Instruct. (a)** Images are encoded by visual embedding layers and then passsed through an abstractor module to reduce token numbers, and then **(c)** fused with text embeddings into under the image-text interleaved format.

### 3.3    Rationale of Combinations.

As discussed in principle 3, our motivation is to collect subsets that can complement each other. This complementarity is reflected in the following two aspects. Firstly, in terms of general comparisons, **Merge2Compare** has higher accuracy but lacks fine-grained comparison (excluded by *Top-similarity Pair Removal*), while **Teach2Compare**-*general*, although slightly less accurate, offers more diverse scenarios and includes content information as background. Joint training of both contributes to a more comprehensive quality comparison by our model. Additionally, **Teach2Compare** includes a unique *Q&A* subset, which significantly enhances the model's ability to answer open-range questions.

## 4    The Co-Instruct Model

In this section, we discuss the proposed model, **Co-Instruct**. Specifically, we have made two non-trivial adaptations for the multi-image comparative setting:

*Visual Token Reduction (Fig. 6 (a)).* Most state-of-the-art LMMs [5,25,39] have adopted the simple projector that keeps a large number of tokens (*e.g.* 1,025). This structure is not friendly to multi-image scenarios: passing only two images will exceed the max length of LLaVA [26] (2,048), and four images will exceed the context window of LLaMA-2 [42] (4,096). Thus, we adopt another widely-used abstractor [4,55,56] structure to reduce token numbers before feeding the visual embeddings to LLM, so as to easily adapt to multi-image scenarios.

*Image-text Interleaved Format (Fig. 6 (c)).* Typical single-image instruction tuning usually does not care about "*position of images*". Most approaches [4,25,68] directly pile all images before texts (`<img`$_0$`>`(`<img`$_1$`>...`)`<text>`). Under this piling, multiple images are not separates and LMMs might confuse the information from different images and fail to compare well (see baseline result in Fig. 8). To solve this, we propose an image-text interleaved format for multi-image training, that each image is started with explicit text to identify its nominal:

`User: The first image:` `<img`$_0$`>` `The second image:` `<img`$_1$`>` `(...)` `<query>`
`Assistant: <response>`

In our experiments, we demonstrated that this interleaved format significantly enhances the performance of **Co-Instruct** (Tab. 8), notably better than using learnable special tokens (`<img_st>` and `<img_end>`) to divide images.
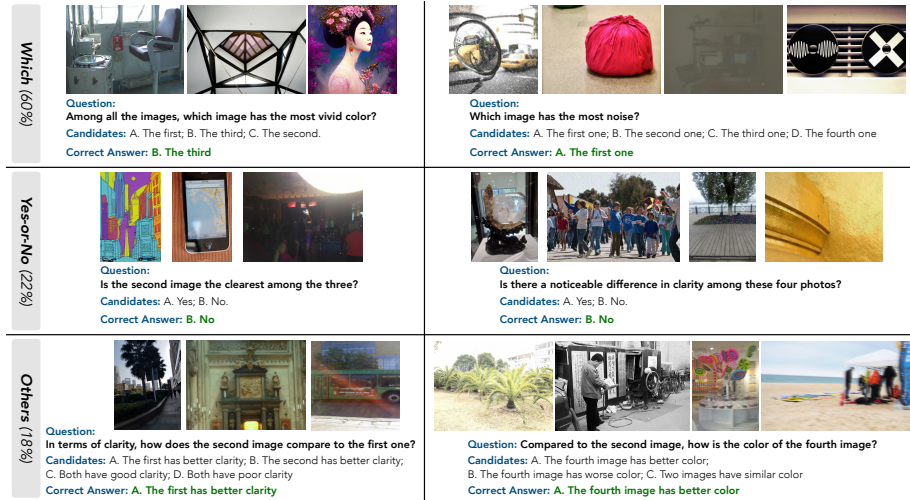
**Fig. 7: Dataset Card of MICBench**, made up of (a) *Which* questions (60%), (b) *Yes-or-No* questions (22%), and (c) *Other* types of questions (18%) on three/four images.

## 5   The MICBench

In this section, we discuss the proposed **MICBench** to cover the open-ended evaluation settings on groups of three or four images, as a complementary of existing evaluation settings (Sec. 6.3). It contains 2,000 groups of *open-range* questions equipped with multiple candidates, with details elaborated as follows:

*Sourcing Diverse Image Groups.* To improve the diversity of the benchmark, in the **MICBench**, we sample image groups from two sources: **(1)** 400 groups of three and 400 groups of four from the images in LLVisionQA [49], which are originally sourced from 9 datasets [3, 6, 7, 12, 15, 18, 24, 59, 66]; **(2)** 600 groups of three and 600 groups of four on 1,000 random images sampled from unlabeled databases [3, 24, 41, 53] (*zero overlap with training-set images*). With in-total 2,000 groups, the **MICBench** contains a wide variety of quality issues and low-level appearances, providing a non-biased evaluation on quality comparison.

*Evaluation Form: Multi-choice Questions (MCQs).* As the most popular evaluation form for LLM/LMM benchmarks [11, 29, 49, 60], multi-choice question (MCQ) is adopted as the evaluation form of **MICBench**. As is shown in Fig. 7, each image group is associated with a expert-crafted question that compare quality or related attributes among the images. Despite common question types (*Yes-or-No/What/How, etc*), the **MICBench** also introduces a special type of question, the ***Which*** questions (Fig. 7(a)), to cover this common type of human query on comparison. In total 10 human experts participate in annotating the **MICBench**, and the answer of each MCQ is cross-examined by another expert. Similar as existing benchmarks [29, 49], **MICBench** is further divided into a *dev* set (1,004) for method development (*answers will be public*), and a *test* set (996) to evaluate performance of LMMs (*answers will be hidden from public*).

## 6      Evaluation

### 6.1      Implementation Details

The **Co-Instruct** is fine-tuned after the released checkpoint of mPLUG-Owl2 [56], with LLaMA-2 [42] as LLM and CLIP-ViT-L14 [36] as visual embedding module. Images are padded to square and then resized to $448 \times 448$ before fed into the model. The learning rate is set as $2e$-5, with two epochs under batch size 192. The final checkpoint is used for evaluation. To avoid over-fitting, only *dev* subsets of evaluation datasets are used to choose best training hyper-parameters, where the final reported results are from non-overlapped *test* subsets. All parameters are updated during training, costing in total 25 hours on 8*NVIDIA A100 GPUs.

### 6.2      Baseline Models

We choose five open-source recent state-of-the-art LMMs that supports multi-image inputs to compare with: LLaVA-v1.5-13B [25], InternLM-XComposer2 [5], BakLLaVA [39], EMU2-Chat [40], mPLUG-Owl2 [56] (*baseline of* **Co-Instruct**). Additionally, we also compare with three well-recognized proprietary close-source models: Qwen-VL-Max, Gemini-Pro, and GPT-4V (*teacher of* **Co-Instruct**).

### 6.3      Results on Existing Evaluation Settings

Despite the **MICBench** (Sec. 5), we also evaluate the proposed **Co-Instruct** against baseline models on several existing visual quality evaluation/comparison benchmarks for LMMs. The evaluation settings and results are as follows.

**Table 2:** Results on *Q-Bench*<sup>PAIR</sup>*-A1*. **Co-Instruct** is remarkably **51%** better than the variant *without comparative data*, and the only LMM that surpasses human capability.
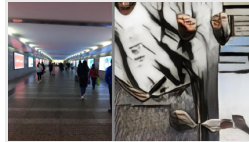
| Sub-categories | Question Types | | | Low-level Concerns | | Pairwise Settings | | |
|---|---|---|---|---|---|---|---|---|
| Model | *Yes-or-No↑* | *What↑* | *How↑* | *Distortion↑* | *Other↑* | *Compare↑* | *Joint↑* | *Overall↑* |
| *random guess accuracy* | 50.00% | 32.03% | 33.16% | 38.95% | 41.95% | 38.69% | 43.70% | 39.82% |
| (Sep/2023) LLaVA-v1.5-13B | 57.34% | 47.45% | 49.13% | 49.01% | 59.51% | 52.06% | 52.00% | 52.05% |
| (Oct/2023) BakLLava | 60.09% | 45.42% | 50.86% | 53.09% | 58.82% | 54.52% | 55.55% | 52.75% |
| (Nov/2023) mPLUG-Owl2 *(baseline of* **Co-Instruct***)* | 58.07% | 36.61% | 48.44% | 47.74% | 51.90% | 45.73% | 60.00% | 48.94% |
| (Dec/2023) Emu2-Chat | 51.94% | 29.78% | 53.84% | 42.01% | 55.71% | 46.26% | 49.09% | 47.08% |
| (Feb/2024) InternLM-XComposer2-VL | 71.81% | 58.64% | 62.28% | 65.77% | 63.67% | 64.34% | **68.00%** | 65.16% |
| Qwen-VL-Max *(Proprietary)* | 67.65% | 67.56% | 65.35% | 69.09% | 61.18% | 68.65% | 61.29% | 66.99% |
| Gemini-Pro *(Proprietary)* | 65.78% | 56.61% | 56.74% | 60.42% | 60.55% | 60.46% | 60.44% | 60.46% |
| GPT-4V *(Proprietary, teacher of* **Co-Instruct***)* | **79.75%** | 69.49% | **84.42%** | 77.32% | 79.93% | 81.00% | 68.00% | 78.07% |
| *Non-expert Human* | 78.11% | 77.04% | 82.33% | 78.17% | 77.22% | 80.26% | 76.39% | 80.12% |
| *without Multi-image Comparative Data* | 60.24% | 47.46% | 48.78% | 52.81% | 53.97% | 51.42% | 59.11% | 53.15% |
| **Co-Instruct** (Ours) | **86.50%** | **72.20%** | 79.23% | **80.00%** | **80.62%** | **81.91%** | 74.22% | **80.18%** |

**Q-Bench**<sup>PAIR</sup>**-A1** [67] is a benchmark for visual quality comparison with 1,999 expert-crafted *open-range* quality-related MCQs on *image pairs*. In Tab. 2, we compare **Co-Instruct** against existing open-source and proprietary models on this benchmark. **Co-Instruct** shows far superior accuracy than open-source LMMs: it is **64%** better than its baseline (mPLUG-Owl2), **51%** better than the variant without our multi-image subsets (**Merge2Compare** and **Teach2Compare**), and also 23% better than the best of them. It also outperforms Qwen-VL-Max

and Gemini-Pro by a large margin (21%/33%). Additionally, though its all MCQ training data are from GPT-4V, the student (**Co-Instruct**) still outperforms its teacher on this MCQ evaluation set by notable **2.7%**, suggesting the effectiveness of the collaborative teaching strategy. Our model is also *the only LMM* that surpasses the accuracy of a non-expert human (*esp.* on Compare subset) in this benchmark, strongly supporting the meaningful vision of using models to relieve human labors on real-world visual quality comparisons in the future.

**Table 3:** Results on $Q\text{-}Bench^{\texttt{PAIR}}\text{-}A2$. $P_i$ denotes frequency for score $i$ (score in $[0, 2]$). While slightly inferior to GPT-4V, **Co-Instruct** has significantly improved over both the variant *without comparative data* (+**31%**), especially for **Precision** metric (+**59%**).

| Dimensions | Completeness | | | | Precision | | | | Relevance | | | | Sum.↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | $P_0$ | $P_1$ | $P_2$ | score↑ | $P_0$ | $P_1$ | $P_2$ | score↑ | $P_0$ | $P_1$ | $P_2$ | score↑ | |
| (Sep/2023) LLaVA-v1.5-13B | 18.77% | 73.44% | 7.79% | 0.89 | 34.66% | 38.72% | 26.62% | 0.92 | 1.02% | 34.59% | 64.39% | 1.63 | 3.44 |
| (Oct/2023) BakLLava | 29.46% | 59.77% | 10.57% | 0.80 | 40.0% | 38.08% | 21.33% | 0.80 | 2.26% | 15.06% | 82.04% | 1.79 | 3.40 |
| (Nov/2023) mPLUG-Owl2 *(baseline)* | 19.43% | 65.54% | 14.45% | 0.94 | 30.94% | 43.71% | 24.63% | 0.92 | 3.79% | 26.94% | 68.28% | 1.63 | 3.50 |
| (Dec/2023) Emu2-Chat | 41.25% | 54.33% | 4.42% | 0.63 | 38.11% | 36.41% | 25.48% | 0.87 | 4.12% | 38.61% | 57.27% | 1.53 | 3.03 |
| (Feb/2024) InternLM-XComposer2-VL | 13.20% | 72.17% | 14.13% | 1.00 | 31.28% | 42.13% | 25.77% | 0.93 | 1.60% | 24.17% | 72.93% | 1.70 | 3.64 |
| Qwen-VL-Max *(Proprietary)* | 11.64% | 54.08% | 34.08% | 1.22 | 24.26% | 39.15% | 36.22% | 1.11 | 2.533% | 10.97% | 85.64% | 1.82 | 4.16 |
| Gemini-Pro *(Proprietary)* | 18.22% | 44.48% | 36.84% | 1.18 | 34.13% | 37.95% | 27.02% | 0.92 | 0.67% | 5.91% | 92.22% | 1.90 | 4.00 |
| GPT-4V *(Proprietary, teacher of **Ours**)* | 4.09% | 31.82% | 64.09% | **1.60** | 10.44% | 45.12% | 44.44% | **1.34** | 0.18% | 1.69% | 96.35% | **1.94** | **4.89** |
| *w/o Multi-Image Comparative Data* | 15.25% | 65.76% | 18.32% | 1.02 | 39.44% | 40.18% | 19.62% | 0.79 | 0.09% | 9.86% | 89.02% | 1.87 | 3.69 |
| **Co-Instruct** (Ours) | 4.04% | 31.55% | 63.55% | **1.58** | 13.68% | 43.68% | 41.37% | **1.26** | 0.0% | 0.44% | 98.22% | **1.96** | **4.82** |



**Fig. 8: Qualitative Visualization on** $Q\text{-}Bench^{\texttt{PAIR}}\text{-}A2$. GPT-4V gives longest outputs and achieves high precision score even if it includes incorrect information.

**Q-Bench$^{\texttt{PAIR}}$-A2** is a benchmark setting for general and detailed visual quality comparison *with detailed reasonings* on *image pairs*. Consisting of 499 image pairs, it employs GPT to evaluate LMM responses against the *golden* expert-labeled comparisons on **Completeness**, **Precision**, and **Relevance**. As listed in Tab. 3, the **Co-Instruct** can remarkably improve the **Completeness** (+57%) and **Precision** (+59%) of the comparative outputs than the *w/o comparative data* version, but still falls a little bit behind GPT-4V on the quantitative metrics. This might be because outputs from GPT-4V are more than *twice as long* as **Co-Instruct** outputs, while GPT evaluation used here is observed [22] to be in favor of longer text outputs. To further analyze this potential bias, we qualitatively visualize the result of different LMMs in Fig. 8. As shown in the figure, the

baseline open-source LMM even confuses the information from the two images, and Gemini-Pro makes rather poor detailed reasonings. For GPT-4V, it generates the *__longest outputs__* among all LMMs, which might be the reason that it gets a relatively high precision score even its outputs are not totally correct. In short, the capability of **Co-Instruct** in reasoning-related comparisons can match that of GPT-4V, while significantly surpassing other existing LMMs.

**Table 4:** Results on *2AFC-LMM*. $\kappa$ denotes binary judgment consistency while swapping *first image* and *second image*; $\rho$ denotes Pearson's linear correlation.

| Dataset | CSIQ | | MM21 | | KADID-10k | | LIVEC | | KonIQ-10k | | SPAQ | | weighted avg. | |
| Model | $\kappa$ | $\rho$ | $\kappa$ | $\rho$ | $\kappa$ | $\rho$ | $\kappa$ | $\rho$ | $\kappa$ | $\rho$ | $\kappa$ | $\rho$ | $\kappa$ | $\rho$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| (Aug/2023) IDEFICS-Instruct-9B | 0.206 | 0.570 | 0.337 | 0.338 | 0.202 | 0.552 | 0.323 | 0.492 | 0.251 | 0.479 | 0.330 | 0.474 | 0.286 | 0.470 |
| (Sep/2023) LLaVA-v1.5-13B | 0.483 | 0.423 | 0.356 | 0.149 | 0.310 | 0.137 | 0.273 | 0.162 | 0.262 | 0.403 | 0.291 | 0.156 | 0.302 | 0.224 |
| (Oct/2023) BakLLava | 0.356 | 0.235 | 0.337 | 0.244 | 0.245 | 0.166 | 0.296 | 0.159 | 0.185 | 0.217 | 0.274 | 0.146 | 0.261 | 0.185 |
| (Nov/2023) mPLUG-Owl2 *(baseline)* | 0.435 | 0.627 | 0.378 | 0.306 | 0.402 | 0.443 | 0.375 | 0.441 | 0.386 | 0.417 | 0.362 | 0.356 | 0.460 | 0.397 |
| (Feb/2024) InternLM-XComposer2-VL | 0.800 | 0.527 | 0.688 | 0.377 | 0.600 | 0.552 | 0.600 | 0.516 | 0.825 | 0.581 | 0.700 | 0.755 | 0.705 | 0.567 |
| Qwen-VL-Max *(Proprietary)* | 0.540 | 0.418 | 0.497 | 0.304 | 0.625 | 0.406 | 0.578 | 0.544 | 0.631 | 0.610 | 0.592 | 0.718 | 0.592 | 0.540 |
| Gemini-Pro *(Proprietary)* | 0.672 | 0.527 | 0.604 | 0.377 | 0.650 | 0.552 | 0.650 | 0.516 | 0.652 | 0.581 | 0.671 | 0.755 | 0.678 | 0.622 |
| GPT-4V *(Proprietary, teacher of **Ours**)* | 0.778 | 0.764 | 0.792 | **0.474** | 0.763 | 0.560 | 0.837 | 0.685 | 0.835 | 0.800 | 0.871 | 0.876 | 0.823 | 0.721 |
| *w/o Multi-Image Comparative Data* | 0.117 | 0.650 | 0.480 | 0.392 | 0.397 | 0.466 | 0.327 | 0.432 | 0.489 | 0.512 | 0.485 | 0.397 | 0.432 | 0.449 |
| **Co-Instruct** (Ours) | **0.800** | **0.779** | **0.852** | 0.325 | **0.829** | **0.685** | **0.872** | **0.797** | **0.883** | **0.927** | **0.881** | **0.931** | **0.864** | **0.754** |

**2AFC-LMM** [69] is a benchmark setting for general quality comparison on *image pairs*. It prompts LMMs to make a *two-alternative forced choice* (2AFC) on a pair of images. The maximum a posterior estimation is utilized to aggregate comparative preferences to single-image quality scores [43]. Then, it computes Peason's linear correlation ($\rho$) between regressed scores and ground truth MOS. As shown in Tab. 4, **Co-Instruct** outperforms all existing models in *2AFC-LMM*, including GPT-4V. **Co-Instruct** also shows very high consistency $\kappa$ while swapping two images. Among all datasets, the proposed model is only inferior on the MM21 [23]. Nonetheless, we observe that the **Co-Instruct** has higher direct comparison accuracy than GPT-4V (**Co-Instruct**: 55.2%, GPT-4V: 54.4%, see supplementary for full results) on it, yet the dataset contains a large proportion of *extremely similar pairs*, for which **Co-Instruct** responds a forced choice but GPT-4V will answer *"two images have similar quality"* (tie), which impacts the aggregated single-image quality scores. We hope this observation can inspire further research to design better evaluation settings for fine-grained comparisons.

**Table 5:** Results on *Q-Bench*<sup>SINGLE</sup>*-A1*, proving that the comparative data (Sec. 3) can also effectively boost the capability of LMMs on single image quality evaluation.

| Sub-categories | Question Types | | | Quadrants of Low-level Concerns | | | | |
| Model | Yes-or-No↑ | What↑ | How↑ | Distortion↑ | Other↑ | In-context Distortion↑ | In-context Other↑ | Overall↑ |
|---|---|---|---|---|---|---|---|---|
| *random guess accuracy* | 50.00% | 28.48% | 33.30% | 37.24% | 38.50% | 39.13% | 37.10% | 37.94% |
| (Sep/2023) LLaVA-v1.5-13B | 64.96% | 64.86% | 54.12% | 53.55% | 66.59% | 58.90% | 71.48% | 61.40% |
| (Oct/2023) BakLLava | 66.46% | 61.48% | 54.83% | 51.33% | 63.76% | 56.52% | 78.16% | 61.02% |
| (Nov/2023) mPLUG-Owl2 *(baseline of **Co-Instruct**)* | 72.26% | 55.53% | 58.64% | 52.59% | 71.36% | 58.90% | 73.00% | 62.68% |
| (Dec/2023) Emu2-Chat | 70.09% | 65.12% | 54.11% | 66.22% | 62.96% | 63.47% | 73.21% | 64.32% |
| (Feb/2024) InternLM-XComposer2-VL | 72.44% | 78.13% | 67.28% | 68.00% | **75.65%** | 68.15% | 81.36% | 72.52% |
| Qwen-VL-Max *(Proprietary)* | 73.20% | **81.02%** | 68.39% | 70.84% | 74.57% | 73.11% | 80.44% | 73.90% |
| Gemini-Pro *(Proprietary)* | 71.26% | 71.39% | 65.59% | 67.30% | 73.04% | 65.88% | 73.60% | 69.46% |
| GPT-4V *(Proprietary, teacher of **Co-Instruct**)* | 77.72% | 78.39% | 66.45% | **71.01%** | 71.07% | **79.36%** | 78.91% | **74.10%** |
| *Non-expert Human* | 82.48% | 79.39% | 60.29% | 75.62% | 72.08% | 76.37% | 73.00% | 74.31% |
| *without Multi-image Comparative Data* | **79.38%** | 72.23% | 67.70% | 68.71% | 72.32% | 73.97% | **83.65%** | 73.38% |
| **Co-Instruct** (Ours) | **81.93%** | 78.74% | **70.16%** | **74.28%** | 76.37% | 76.71% | **84.41%** | **77.12%** |

**Q-Bench**<sup>SINGLE</sup>**-A1.** Despite the comparative benchmarks above, we also evaluate **Co-Instruct** on *single image* MCQs from *Q-Bench*<sup>SINGLE</sup>*-A1* to verify the

influences of comparative tuning on single-image quality perception. As shown in Tab. 5, **Co-Instruct** shows **5%** improvement over the variant *trained with single images only*, leading GPT-4V by 4%, and marks the only LLM that surpasses non-expert human. These results have demonstrated the contribution of comparative training on general quality-related understanding of LMMs, and suggested that single-image quality evaluation *does not conflict with* multi-image quality comparison for LMMs and can be improved together under a unified model.

## 6.4    Results on MICBench

**Table 6:** Results on **MICBench**. **Co-Instruct** is **60%** better than the variant *without comparative data*, and also notably better than GPT-4V ($+5.7\%$) and human ($+6.4\%$).

| Sub-categories | Question Types | | | Number of Images | | Overall↑ |
|---|---|---|---|---|---|---|
| Model | *Yes-or-No*↑ | *Which*↑ | *Others*↑ | *Three*↑ | *Four*↑ | |
| #questions | 220 | 594 | 182 | 503 | 493 | 996 |
| *random guess accuracy* | 49.55% | 28.59% | 28.31% | 34.10% | 29.17% | 31.47% |
| (Sep/2023) LLaVA-v1.5-13B (*length: 2048→2560*) | 47.51% | 40.74% | 52.49% | 46.81% | 41.90%* | 44.38% |
| (Oct/2023) BakLLava (*length: 2048→2560*) | 68.35% | 35.01% | 52.78% | 48.51% | 42.54%* | 45.56% |
| (Nov/2023) mPLUG-Owl2 (*baseline of **Co-Instruct***) | 62.25% | 35.70% | 53.71% | 44.19% | 45.42% | 44.80% |
| (Feb/2024) InternLM-XComposer2-VL (*length: 4096→5120*) | 62.95% | 47.29% | 52.02% | 55.70% | 46.51%* | 51.76% |
| Qwen-VL-Max *(Proprietary)* | 62.33% | 70.00% | **81.54%** | 72.35% | 68.79% | 70.55% |
| Gemini-Pro *(Proprietary)* | 75.00% | 67.37% | 66.92% | 68.71% | 70.87% | 69.79% |
| GPT-4V *(Proprietary, teacher of **Co-Instruct**)* | **80.32%** | **77.28%** | 78.82% | **80.32%** | **77.28%** | **78.82%** |
| *Non-expert Human* | 82.27% | 78.15% | 74.31% | 77.18% | 79.55% | 78.35% |
| *without Multi-image Comparative Data* | 62.72% | 37.54% | 53.30% | 45.33% | 46.65% | 45.98% |
| **Co-Instruct** (Ours) | 79.55% | **85.35%** | **81.32%** | **84.69%** | **81.94%** | **83.33%** |

As is shown in Tab. 6, **Co-Instruct** provides very competitive accuracy on open-question quality comparison among three/four images, **5.7%** better than GPT-4V (*best existing*) and **6.4%** more accurate than non-expert human; open-source LMMs even struggle to obtain 50% accuracy on this setting. It is also noteworthy that LLaVA series and XComposer2-VL's original context lengths are *not enough for four images* as they have not reduced visual token numbers, so we have to extend their context windows to evaluate them for **MICBench**. Consequentially, all these models have experienced notably worse accuracy on groups of four (*on extended context length*) than groups of three (*within original context length*), as noted in * in Tab. 6. This degradation highlights the importance to *reduce visual tokens* (Sec. 4) to adapt to multi-image scenarios.
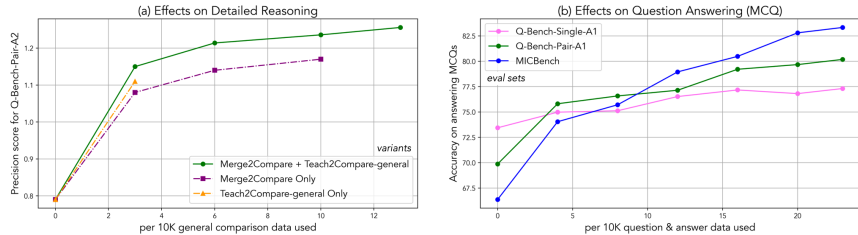
## 6.5    Ablation Studies

**Ablation on Training Data.** As the proposed dataset is composed of three parts, we discuss the effects on different subsets of data on the six evaluation scenarios in Tab. 7. Through the results, we have arrived at several important conclusions: **A)** Even with all information from Q-Instruct-200K, the incorporation of **Merge2Compare** (variant (2)) still notably enhances the capability across various settings; **B)** Only involving the **Teach2Compare** (*i.e.* GPT-4V labels, variant (5)) cannot outperform its teacher GPT-4V; **C)** instead, the superiority towards GPT-4V on question-answering benchmarks is benefited from the co-instruction on *hetero-sourced* training subsets (*variants (6) and (7)*).

**Table 7: Ablation Study** on data. *Abbreviation:* Merge: **Merge2Compare**; Teach$^G$: **Teach2Compare**-*general*; Teach$^{QA}$: **Teach2Compare**-*Q&A*; QIn: Q-Instruct-200K.

| Training Subset | | | | | Evaluation Scenario | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Variant No. | QIn | Merge | Teach$^G$ | Teach$^{QA}$ | Q-Bench$^{\text{PAIR}}$-A1 | Q-Bench$^{\text{PAIR}}$-A2 | Q-Bench$^{\text{SINGLE}}$-A1 | 2AFC-LMM | **MICBench** |
| *Reference Results of GPT-4V* | | | | | 78.07 | 4.89 | 74.10 | 0.721 | 78.82 |
| (1) | ✓ | | | | 53.15 | 3.69 | 73.38 | 0.449 | 45.98 |
| (2) | ✓ | ✓ | | | 68.67 | 4.67 | 75.12 | 0.701 | 60.34 |
| (3) | ✓ | | ✓ | | 65.44 | 4.64 | 74.38 | 0.647 | 64.82 |
| (4) | ✓ | ✓ | ✓ | | 69.87 | **4.82** | 76.52 | 0.749 | 66.37 |
| (5) | ✓ | | ✓ | ✓ | 78.28 | 4.65 | 75.72 | 0.676 | 76.41 |
| (6) | ✓ | ✓ | | ✓ | 80.08 | 4.68 | 75.65 | 0.728 | 81.82 |
| (7, *full*) | ✓ | ✓ | ✓ | ✓ | **80.18** | **4.82** | **77.31** | **0.754** | **83.33** |

**Table 8: Ablation Study** on the text-image interleaved format for the **Co-Instruct**.

| Format | Q-Bench$^{\text{PAIR}}$-A1 | Q-Bench$^{\text{PAIR}}$-A2 | Q-Bench$^{\text{SINGLE}}$-A1 | 2AFC-LMM | **MICBench** |
|---|---|---|---|---|---|
| `<img₀><img₁> ...)` (*baseline*, popular strategy [4,25,68]) | 76.37 | 4.73 | 74.11 | 0.729 | 78.92 |
| `<img_st><img₀><img_end> (<img_st><img₁><img_end>...)` | 76.77 | 4.79 | 73.85 | 0.736 | 80.12 |
| `The input image: <img₀> (The input image: <img₁>...)` | 78.28 | 4.80 | 75.12 | 0.749 | 82.22 |
| `The first image: <img₀> (The second image: <img₁>...)` | **80.18** | **4.82** | **77.31** | **0.754** | **83.33** |



**Fig. 9: Effects of Data Scaling** for (a) *detailed reasoning* and (b) *question answering.*

Besides the composition of three subsets, we also explore the effects of data scale: in Fig. 9(a), we confirm that more general comparison data contributes to an increase in detailed reasoning capabilities, and mixed data has a better effect than homogenous data at the same scale; in Fig. 9(b), we also validate that scaling up the Q&A subset helps to improve multiple MCQ metrics.

**Ablation on Interleaved Format.** As shown in Tab. 8, the proposed text-image interleaved format has proved non-negligible advantages than the baseline, as well as other variants *without explicitly noting* image orders. The results suggest the rationale of the format on multi-image comparative settings.

## 7    Conclusion

In this work, we investigate the *open-ended visual quality comparison* problem, with the aim of a model that provides answers and *detailed reasonings* on *open-range questions* that compares quality among multiple images. To achieve this, we collect the first instruction-tuning dataset to fine-tune large multi-modality models (LMMs) for comparison, the Co-Instruct-562K, with two subsets from human annotations on single images (*merged by LLMs*), and GPT-4V responses. With the dataset, we propose the **Co-Instruct**, which not only outperforms all existing LMMs (including *its teacher*, GPT-4V) on visual quality comparison, but also marks the first LMM with the capability to surpass human accuracy on related settings. We further construct the **MICBench**, the first benchmark that evaluates multi-image quality comparison for LMMs on three and four images. We expect our work to motivate future studies on visual quality comparison.

# References

1. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: VQA: Visual Question Answering. In: IEEE ICCV. pp. 2425–2433 (2015) 2

2. Chen, L., Li, J., Dong, X., Zhang, P., He, C., Wang, J., Zhao, F., Lin, D.: ShareGPT4V: Improving large multi-modal models with better captions. CoRR **abs/2311.12793** (2023) 3

3. Chen, X., Fang, H., Lin, T.Y., Vedantam, R., Gupta, S., Dollár, P., Zitnick, C.L.: Microsoft COCO captions: Data collection and evaluation server. CoRR **abs/1504.00325** (2015) 7, 9

4. Dai, W., Li, J., Li, D., Tiong, A.M.H., Zhao, J., Wang, W., Li, B., Fung, P., Hoi, S.: InstructBLIP: Towards general-purpose vision-language models with instruction tuning. CoRR **abs/2305.06500** (2023) 2, 8, 14

5. Dong, X., Zhang, P., Zang, Y., Cao, Y., Wang, B., Ouyang, L., Wei, X., Zhang, S., Duan, H., Cao, M., Zhang, W., Li, Y., Yan, H., Gao, Y., Zhang, X., Li, W., Li, J., Chen, K., He, C., Zhang, X., Qiao, Y., Lin, D., Wang, J.: InternLM-XComposer2: Mastering free-form text-image composition and comprehension in vision-language large model. CoRR **abs/2401.16420** (2024) 2, 3, 8, 10

6. Fang, Y., Zhu, H., Zeng, Y., Ma, K., Wang, Z.: Perceptual quality assessment of smartphone photography. In: IEEE CVPR. pp. 3677–3686 (2020) 9

7. Ghadiyaram, D., Bovik, A.C.: Massive online crowdsourced study of subjective and objective picture quality. IEEE TIP **25**(1), 372–387 (2016) 9

8. Golestaneh, S.A., Dadsetan, S., Kitani, K.M.: No-reference image quality assessment via transformers, relative ranking, and self-consistency. In: IEEE WACV. pp. 3209–3218 (2022) 4

9. Google: Gemini Pro (2023), https://deepmind.google/technologies/gemini 2

10. Gu, J., Cai, H., Chen, H., Ye, X., Ren, J., Dong, C.: PIPAL: A large-scale image quality assessment dataset for perceptual image restoration. In: ECCV. pp. 633–651 (2020) 2, 5

11. Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., Steinhardt, J.: Measuring massive multitask language understanding. In: ICLR. pp. 1–10 (2021) 9

12. Hosu, V., Lin, H., Sziranyi, T., Saupe, D.: Koniq-10k: An ecologically valid database for deep learning of blind image quality assessment. IEEE TIP **29**, 4041–4056 (2020) 9

13. Hudson, D.A., Manning, C.D.: GQA: A new dataset for real-world visual reasoning and compositional question answering. In: IEEE CVPR. pp. 6700–6709 (2019) 2

14. ITU-R, B.T.: Methodology for the subjective assessment of the quality of television pictures. https://www.itu.int/rec/R-REC-BT.500 (2002) 4

15. Jayaraman, D., Mittal, A., Moorthy, A.K., Bovik, A.C.: Objective quality assessment of multiply distorted images. In: ASILOMAR. pp. 1693–1697 (2012) 9

16. Jiang, A.Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D.S., de las Casas, D., Hanna, E.B., Bressand, F., Lengyel, G., Bour, G., Lample, G., Lavaud, L.R., Saulnier, L., Lachaux, M.A., Stock, P., Subramanian, S., Yang, S., Antoniak, S., Scao, T.L., Gervet, T., Lavril, T., Wang, T., Lacroix, T., Sayed, W.E.: Mixtral of experts. CoRR **abs/2401.04088** (2024) 3, 6

17. LAION: LAION GPT-4V dataset. https://huggingface.co/datasets/laion/gpt4v-dataset (2023) 3

18. Li, C., Zhang, Z., Wu, H., Sun, W., Min, X., Liu, X., Zhai, G., Lin, W.: AGIQA-3K: An open database for ai-generated image quality assessment. CoRR **2306.04717** (2023) 9

19. Li, D., Jiang, T., Jiang, M.: Quality assessment of in-the-wild videos. In: ACM MM. pp. 2351–2359 (2019) 7
20. Li, D., Jiang, T., Lin, W., Jiang, M.: Which has better visual quality: The clear blue sky or a blurry animal? IEEE TMM **21**(5), 1221–1234 (2019) 7
21. Li, J., Mantiuk, R., Wang, J., Ling, S., Le Callet, P.: Hybrid-MST: A hybrid active sampling strategy for pairwise preference aggregation. In: NeurIPS. pp. 1–11 (2018) 4, 5
22. Li, X., Zhang, T., Dubois, Y., Taori, R., Gulrajani, I., Guestrin, C., Liang, P., Hashimoto, T.B.: AlpacaEval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval (2023) 11
23. Li, Y., Wang, S., Zhang, X., Wang, S., Ma, S., Wang, Y.: Quality assessment of end-to-end learned image compression: The benchmark and objective measure. In: ACM MM. pp. 4297–4305 (2021) 12
24. Lin, H., Hosu, V., Saupe, D.: KADID-10k: A large-scale artificially distorted iqa database. In: QoMEX. pp. 1–3 (2019) 7, 9
25. Liu, H., Li, C., Li, Y., Lee, Y.J.: Improved baselines with visual instruction tuning. CoRR **abs/2310.03744** (2023) 2, 3, 8, 10, 14
26. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning. CoRR **abs/2304.08485** (2023) 2, 3, 6, 8
27. Liu, X., Van De Weijer, J., Bagdanov, A.D.: RankIQA: Learning from rankings for no-reference image quality assessment. In: IEEE ICCV. pp. 1040–1049 (2017) 4
28. Liu, X., Van De Weijer, J., Bagdanov, A.D.: Exploiting unlabeled data in cnns by self-supervised learning to rank. IEEE TPAMI **41**(8), 1862–1878 (2019) 2
29. Liu, Y., Duan, H., Zhang, Y., Li, B., Zhang, S., Zhao, W., Yuan, Y., Wang, J., He, C., Liu, Z., Chen, K., Lin, D.: MMBench: Is your multi-modal model an all-around player? CoRR **abs/2307.06281** (2023) 9
30. Ma, K., Liu, W., Liu, T., Wang, Z., Tao, D.: dipIQ: Blind image quality assessment by learning-to-rank discriminable image pairs. IEEE TIP **26**(8), 3951–3964 (2017) 4
31. Mantiuk, R.K., Tomaszewska, A., Mantiuk, R.: Comparison of four subjective methods for image quality assessment. In: Computer Graphics Forum. vol. 31, pp. 2478–2491 (2012) 4
32. Michaelis, C., Mitzkus, B., Geirhos, R., Rusak, E., Bringmann, O., Ecker, A.S., Bethge, M., Brendel, W.: Benchmarking robustness in object detection: Autonomous driving when winter is coming. CoRR **abs/1907.07484** (2019) 7
33. Mikhailiuk, A., Wilmot, C., Perez-Ortiz, M., Yue, D., Mantiuk, R.K.: Active sampling for pairwise comparisons via approximate message passing and information gain maximization. In: IEEE ICPR. pp. 2559–2566 (2021) 4
34. OpenAI: Gpt-4 technical report (2023) 2
35. Prashnani, E., Cai, H., Mostofi, Y., Sen, P.: PieAPP: Perceptual image-error assessment through pairwise preference. In: IEEE CVPR. pp. 1808–1817 (2018) 2, 5
36. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: ICML. pp. 8748–8763 (2021) 10
37. Rajkumar, A., Agarwal, S.: When can we rank well from comparisons of o (nlog (n)) non-actively chosen pairs? In: Conference on Learning Theory. pp. 1376–1401 (2016) 6
38. Schwenk, D., Khandelwal, A., Clark, C., Marino, K., Mottaghi, R.: A-OKVQA: A benchmark for visual question answering using world knowledge. In: ECCV. pp. 146–162 (2022) 7

39. SkunkworksAI: BakLLaVA (2024), https://github.com/SkunkworksAI/BakLLaVA 3, 8, 10

40. Sun, Q., Cui, Y., Zhang, X., Zhang, F., Yu, Q., Luo, Z., Wang, Y., Rao, Y., Liu, J., Huang, T., et al.: Generative multimodal models are in-context learners. CoRR **abs/2312.13286** (2023) 10

41. Thomee, B., Shamma, D.A., Friedland, G., Elizalde, B., Ni, K., Poland, D., Borth, D., Li, L.J.: YFCC100M: The new data in multimedia research. Commun. ACM **59**(2), 64–73 (2016) 7, 9

42. Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Ferrer, C.C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P.S., Lachaux, M.A., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E.M., Subramanian, R., Tan, X.E., Tang, B., Taylor, R., Williams, A., Kuan, J.X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., Scialom, T.: Llama 2: Open foundation and fine-tuned chat models. CoRR **abs/2307.09288** (2023) 3, 8, 10

43. Tsukida, K., Gupta, M.R.: How to analyze paired comparison data (Technical Report UWEETR-2011-0004, University of Washington, 2011), https://api.semanticscholar.org/CorpusID:15425240 2, 12

44. Wang, L., Yang, N., Huang, X., Yang, L., Majumder, R., Wei, F.: Improving text embeddings with large language models. CoRR **abs/2401.00368** (2024) 3, 6

45. Wauthier, F., Jordan, M., Jojic, N.: Efficient ranking from pairwise comparisons. In: ICML. pp. 109–117 (2013) 6

46. Wu, H., Chen, C., Hou, J., Liao, L., Wang, A., Sun, W., Yan, Q., Lin, W.: FAST-VQA: Efficient end-to-end video quality assessment with fragment sampling. In: ECCV. pp. 538–554 (2022) 4

47. Wu, H., Chen, C., Liao, L., Hou, J., Sun, W., Yan, Q., Gu, J., Lin, W.: Neighbourhood representative sampling for efficient end-to-end video quality assessment. IEEE TPAMI (2023) 7

48. Wu, H., Zhang, E., Liao, L., Chen, C., Hou, J., Wang, A., Sun, W., Yan, Q., Lin, W.: Exploring video quality assessment on user generated contents from aesthetic and technical perspectives. In: IEEE ICCV (2023) 4

49. Wu, H., Zhang, Z., Zhang, E., Chen, C., Liao, L., Wang, A., Li, C., Sun, W., Yan, Q., Zhai, G., Lin, W.: Q-Bench: A benchmark for general-purpose foundation models on low-level vision. In: ICLR (2024) 2, 4, 5, 9

50. Wu, H., Zhang, Z., Zhang, E., Chen, C., Liao, L., Wang, A., Xu, K., Li, C., Hou, J., Zhai, G., et al.: Q-instruct: Improving low-level visual abilities for multi-modality foundation models. CoRR **abs/2311.06783** (2023) 2, 3, 5

51. Wu, H., Zhang, Z., Zhang, W., Chen, C., Liao, L., Li, C., Gao, Y., Wang, A., Zhang, E., Sun, W., et al.: Q-Align: Teaching lmms for visual scoring via discrete text-defined levels. CoRR **abs/2312.17090** (2023) 2, 5

52. Wu, X., Sun, K., Zhu, F., Zhao, R., Li, H.: Better aligning text-to-image models with human preference. CoRR **abs/2303.14420** (2023) 2

53. Xu, J., Liu, X., Wu, Y., Tong, Y., Li, Q., Ding, M., Tang, J., Dong, Y.: ImageReward: Learning and evaluating human preferences for text-to-image generation. CoRR **abs/2304.05977** (2023) 2, 7, 9

54. Ye, P., Doermann, D.: Active sampling for subjective image quality assessment. In: IEEE CVPR. pp. 4249–4256 (2014) 4, 5

55. Ye, Q., Xu, H., Xu, G., Ye, J., Yan, M., Zhou, Y., Wang, J., Hu, A., Shi, P., Shi, Y., Jiang, C., Li, C., Xu, Y., Chen, H., Tian, J., Qi, Q., Zhang, J., Huang, F.: mPLUG-Owl: Modularization empowers large language models with multimodality. CoRR **abs/2304.14178** (2023) 8

56. Ye, Q., Xu, H., Ye, J., Yan, M., Liu, H., Qian, Q., Zhang, J., Huang, F., Zhou, J.: mPLUG-Owl2: Revolutionizing multi-modal large language model with modality collaboration. CoRR **abs/2311.04257** (2023) 2, 3, 8, 10

57. Yim, J.G., Wang, Y., Birkbeck, N., Adsumilli, B.: Subjective quality assessment for youtube UGC dataset. In: IEEE ICIP. pp. 1–5 (2020) 2

58. Yin, Z., Wang, J., Cao, J., Shi, Z., Liu, D., Li, M., Sheng, L., Bai, L., Huang, X., Wang, Z., et al.: LAMM: Language-assisted multi-modal instruction-tuning dataset, framework, and benchmark. CoRR **abs/2306.06687** (2023) 6

59. Ying, Z., Niu, H., Gupta, P., Mahajan, D., Ghadiyaram, D., Bovik, A.: From patches to pictures (PaQ-2-PiQ): Mapping the perceptual space of picture quality. In: IEEE CVPR. pp. 3575–3585 (2020) 9

60. Yue, X., Ni, Y., Zhang, K., Zheng, T., Liu, R., Zhang, G., Stevens, S., Jiang, D., Ren, W., Sun, Y., et al.: MMMU: A massive multi-discipline multimodal understanding and reasoning benchmark for expert AGI. CoRR **abs/2311.16502** (2023) 9

61. Zhang, C., Su, S., Zhu, Y., Yan, Q., Sun, J., Zhang, Y.: Exploring and evaluating image restoration potential in dynamic scenes. In: IEEE CVPR. pp. 2057–2066 (2022) 2

62. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: IEEE CVPR. pp. 586–595 (2018) 2

63. Zhang, W., Ma, K., Zhai, G., Yang, X.: Uncertainty-aware blind image quality assessment in the laboratory and wild. IEEE TIP **30**, 3474–3486 (2021) 2, 4

64. Zhang, W., Zhai, G., Wei, Y., Yang, X., Ma, K.: Blind image quality assessment via vision-language correspondence: A multitask learning perspective. In: IEEE CVPR. pp. 14071–14081 (2023) 4

65. Zhang, W., Liu, Y., Dong, C., Qiao, Y.: RankSRGan: Generative adversarial networks with ranker for image super-resolution. In: ICCV. pp. 3096–3105 (2019) 2

66. Zhang, Z., Sun, W., Wang, T., Lu, W., Zhou, Q., Wang, Q., Min, X., Zhai, G., et al.: Subjective and objective quality assessment for in-the-wild computer graphics images. ACM TOMM **20**(4), 1–22 (2023) 9

67. Zhang, Z., Wu, H., Zhang, E., Zhai, G., Lin, W.: A benchmark for multimodal foundation models on low-level vision: from single images to pairs. CoRR **abs/2402.07116** (2024) 2, 3, 4, 5, 7, 10

68. Zhu, D., Chen, J., Shen, X., Li, X., Elhoseiny, M.: MiniGPT-4: Enhancing vision-language understanding with advanced large language models. CoRR **abs/2304.10592** (2023) 6, 8, 14

69. Zhu, H., Sui, X., Chen, B., Liu, X., Chen, P., Fang, Y., Wang, S.: 2AFC prompting of large multimodal models for image quality assessment. CoRR **abs/2402.01162** (2024) 2, 3, 5, 7, 12