# 🔥FIRE: A Dataset for Feedback Integration and Refinement Evaluation of Multimodal Models

**Pengxiang Li**[1,2*] **Zhi Gao**[2,3*] **Bofei Zhang**[2*] **Tao Yuan**[2] **Yuwei Wu**[1✉]
**Mehrtash Harandi**[4] **Yunde Jia**[1] **Song-Chun Zhu**[2,3,5] **Qing Li**[2✉]
[1]Beijing Institute of Technology [2]Beijing Institute for Artificial General Intelligence
[3]Peking University [4]Monash University [5]Tsinghua University
mm-fire.github.io

Figure 1: The comparison of the **feedback-refining** capability among different models. While the original LLaVA hardly improves its responses, our model trained on FIRE can effectively integrate the user feedback and produce much better responses, which are closer to those of GPT-4V.

## Abstract

Vision language models (VLMs) have achieved impressive progress in diverse applications, becoming a prevalent research direction. In this paper, we build FIRE, a feedback-refinement dataset, consisting of 1.1M multi-turn conversations that are derived from 27 source datasets, empowering VLMs to spontaneously refine their responses based on user feedback across diverse tasks. To scale up the data collection, FIRE is collected in two components: FIRE-100K and FIRE-1M, where FIRE-100K is generated by GPT-4V, and FIRE-1M is freely generated via models trained on FIRE-100K. Then, we build FIRE-Bench, a benchmark to comprehensively evaluate the feedback-refining capability of VLMs, which contains 11K feedback-refinement conversations as the test data, two evaluation settings, and a model to provide feedback for VLMs. We develop the FIRE-LLaVA model by fine-tuning LLaVA on FIRE-100K and FIRE-1M, which shows remarkable feedback-refining capability on FIRE-Bench and outperforms untrained VLMs by 50%, making more efficient user-agent interactions and underscoring the significance of the FIRE dataset.

---

*Equal contribution. ✉ Corresponding author.

# 1 Introduction

Vision language models (VLMs), such as LLaVA [32], GPT-4V [39], and Gemini [46], have demonstrated remarkable progress across various tasks [54, 30, 9] by integrating large language models (LLMs) [48, 20] with visual encoders [12, 42]. However, VLMs can sometimes produce undesirable outputs, possibly due to omitting important details in images or misunderstanding the instructions, which prompts the need for the **feedback-refining** capability beyond the normal instruction-following ability. This capability enables VLMs to spontaneously refine their responses based on user feedback, as depicted in Fig. 1, enhancing the efficiency and smoothness of interactions between users and visual assistants.

In doing so, we build FIRE, a dataset for <u>F</u>eedback <u>I</u>ntegration and <u>R</u>efinement <u>E</u>valuation of VLMs. FIRE is composed of 1.1M high-quality multi-turn feedback-refinement conversations, derived from 27 source datasets across a wide range of tasks, such as visual question answering [15], image captioning [7], OCR reasoning [38, 43], document understanding [17], math reasoning [34], and chart analysis [36]. To scale up the data collection, FIRE is collected in two stages. In the first stage, we randomly sample ∼100K image-instruction-response triplets from data sources. We use each triplet to instruct GPT-4V to simulate a dialogue between a student and a teacher: the student answers the question and the teacher provides feedback to help the student improve its answer. We filter out generated low-quality conversations, such as those with too many turns or no improvement, rendering 100K high-quality feedback-refinement conversations, named FIRE-100K. In the second stage, we fine-tune two LLaVA-NeXT [31] models on FIRE-100K: one is trained as a student to refine its answer with the feedback, and the other is trained as a teacher to generate feedback for the student's answer. We simulate dialogues between the student and the teacher models using ∼1M data points from the sources, rending a split named FIRE-1M. In this case, the full FIRE dataset consists of 1.1M feedback-refinement conversations in two splits FIRE-100K and FIRE-1M.

To comprehensively evaluate the feedback-refining capability of VLMs, we build FIRE-Bench that has 11K feedback-refinement conversations derived from 16 source datasets, including 8 seen datasets (test splits) from FIRE-100K and FIRE-1M, as well as 8 new datasets from recently-proposed popular multimodal benchmarks. Using FIRE-Bench, we design two evaluation settings: fixed dialogues and free dialogues. In fixed dialogues, we compare the model's refined response with ground truth in the generated conversations in FIRE-Bench, given a fixed dialogue history. In free dialogues, we let the model freely interact with a teacher model about instructions in FIRE-Bench, and test how fast & how much the model can improve its answers based on the feedback provided by the teacher model.

We develop FIRE-LLaVA by fine-tuning LLaVA-NeXT [31] on FIRE-100K and FIRE-1M. The evaluation results on FIRE-Bench show that FIRE-LLaVA exhibits significant improvements based on feedback in conversations, exceeding the original LLaVA-Next model method by $50\%$. These results underscore the significance of FIRE-100K and FIRE-1M in enhancing feedback integration, while FIRE-Bench provides an evaluation platform to analyze refinements. We expect that FIRE could motivate future exploration for the feedback-refining capability of VLMs.

In summary, our contributions are three-fold. (1) We introduce FIRE, a dataset containing 1.1M feedback-refinement conversations across a wide range of tasks, where 100K data is generated by GPT-4V and 1M data is freely generated by simulating dialogues between student and teacher models. (2) We introduce the FIRE-Bench benchmark, composed of 11K conversations and a teacher model, providing comprehensive evaluations for the feedback-refining capability in two settings: fixed dialogues and free dialogues. (3) We develop FIRE-LLaVA, an advanced VLM that could improve its responses based on feedback, making efficient interaction between users and VLMs.

# 2 Related Work

## 2.1 Vision Language Models

Building open-source VLMs to compete with closed-source models like GPT-4V [39] and Gemini [46] is a hot research topic. BLIP [24, 23] and Flamingo [1] are pioneering models that combine LLMs with visual encoders to enhance cross-modal understanding and reasoning abilities. LLaVA [32],

InstructBLIP [11], and MiniGPT4 [54] develop the instruction tuning ability of VLMs by introducing a large number of instruction-response pairs. Along this way, some work focuses on the visual grounding ability of VLMs, such as Kosmos-2 [41], MINI-GPTv2 [5], and Qwen-VL [3], improving the region understanding for VLMs. InternVL [9] and mini-Gemini [26] develop powerful visual encoders for high-resolution images, and CuMo adopts a mixture-of-experts (MOE) architecture to better manage diverse data. Compared with existing VLMs, our FIRE-LLaVA has a more powerful feedback-refining capability across diverse tasks, which can spontaneously refine responses based on user feedback, leading to efficient and smooth interaction with users.

## 2.2 Vision-Language Data Generation

Recent attention has increasingly focused on synthesizing vision-language data. The ShareGPT4V dataset [7] leverages GPT-4V to generate 1.2M image-text pairs with detailed descriptions, leading to better alignments. LLaVA-Instruct-150K [32] is a general visual instruction tuning dataset that is constructed by feeding captions and bounding boxes to GPT-4. After that, many efforts have been made to enhance the data diversity of instruction tuning data. LLaVAR [52], MIMIC-IT [22], and SVIT [53] further scale up it to 422K, 2.8M, and 4.2M, respectively. InternLM-XComposer [51] produces interleaved instruction and image data, enabling advanced image-text comprehension and composition. Mini-Gemini [26] and ALLaVA [4] use GPT-4V to exploit visual information and generate high-quality instruction data. LRV-Instruction [29] creates both positive and negative instructions for the hallucinating inconsistent issue. A recent work DRESS [8] collects 66K feedback data and trains VLMs for the feedback-refining capability. Different from DRESS that only uses data from LLaVA-Instruct-150K, our feedback-refinement data is from richer sources (27 datasets) across more tasks (math reasoning, chart understanding, and OCR *etc.*). Moreover, FIRE has significantly more data than DRESS (1.1M *v.s.* 66K), where 1M data is freely produced via dialogues of student and teacher models, leading to significant data expansion but a similar cost of data generation.

## 2.3 Feedback Learning in Multimodal Models

Learning from feedback is a promising research direction, playing an important role in human-robot interaction [27, 10]. Existing feedback learning methods can be roughly divided into two categories: planned feedback learning and impromptu feedback learning. Planned feedback learning updates models based on user feedback, and thus can generalize to new data but cannot provide refined responses immediately. CLOVA [14] and Clarify [21] are representative methods that automatically collect data to learn new knowledge. LLaVA-RLHF [45] collects human preference and trains VLMs via reinforcement learning. Impromptu feedback learning can immediately refine responses but have less generalization since they usually do not update models, which is widely studied in LLMs [2, 25, 47]. Liao *et al.* [28] use VLMs themselves as verifiers that produce feedback to correct recognition results. DRESS [8] generates helpfulness, honesty, and harmlessness responses via impromptu feedback learning. Different from DRESS, we improve the correctness and details of responses via impromptu feedback learning across diverse tasks.
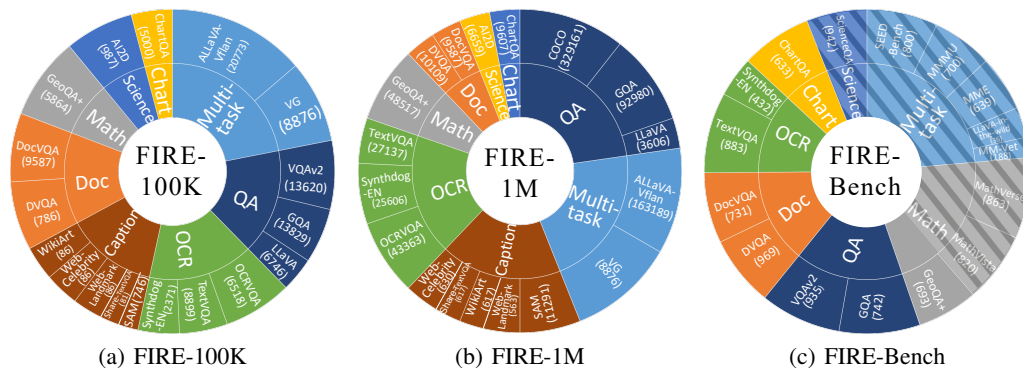


(a) FIRE-100K     (b) FIRE-1M     (c) FIRE-Bench

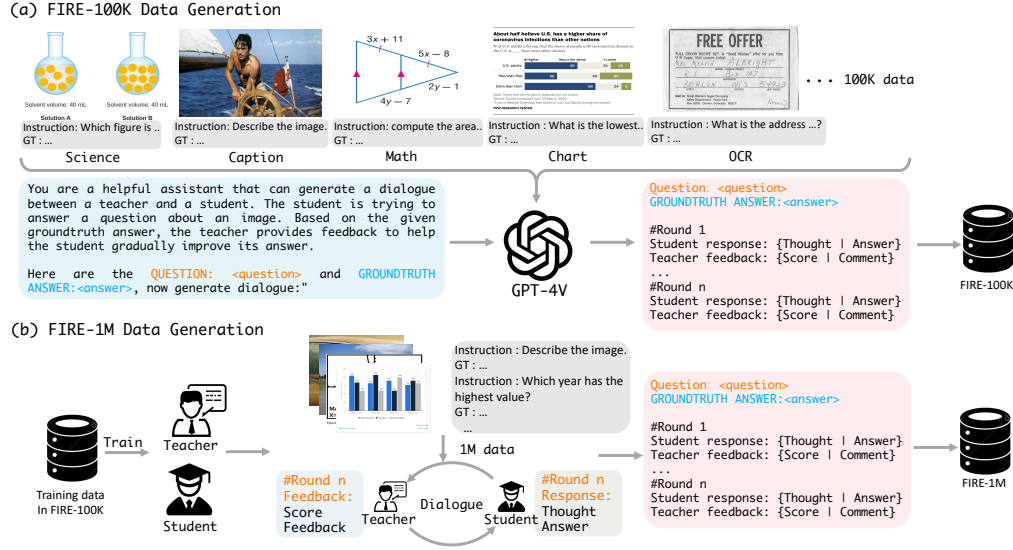Figure 2: Data sources in FIRE. Shaded are new data sources in FIRE-Bench.

Figure 3: The pipeline to create FIRE-100K and FIRE-1M data.

## 3 F̲eedback I̲ntegration and R̲efinement E̲valuation (FIRE)

This section presents the FIRE dataset, outlining its task definition, data collection methodology for FIRE-100K and FIRE-1M, and the creation of FIRE-Bench. Finally, we provide an analysis of FIRE.

### 3.1 Task Definition

**Data Sources** To enhance the diversity and comprehensiveness of our dataset, we compile more than 1.1M image-instruction-response triples from 27 source datasets, being used to generate FIRE-100K, FIRE-1M, and FIRE-Bench, as shown in Fig. 2. These datasets cover tasks including visual question answering, image captioning, complex reasoning, OCR, chart/table/document analysis, math problems, science question answering *etc*.

**Data format.** We formulate our data as $\{I, q, gt, \{r_i, f_i\}_{i=1}^n\}$, where $I$ denotes the image, $q$ is the instruction, $gt$ is the ground truth answer, and $\{r_i, f_i\}_{i=1}^n$ corresponds to the conversations in $n$ turns. In the $i$-th turn, $r_i$ is the response from VLMs, composed of the thought (how to refine the response based on feedback) and a new answer; $f_i$ is the feedback, involving a score $a_i$ (0-10) for the response $r_i$ and textual comments.

### 3.2 FIRE-100K

We feed images, instructions, ground truth answers from $18$ datasets, and a designed textual prompt to GPT-4V that generates high-quality feedback-refinement conversations in a one-go manner, as shown in Fig. 3 (a). We ask GPT-4V to play two roles: a student and a teacher, and generate a conversation between the two roles, where the student's responses are improved by incorporating feedback from the teacher. After generation, we filter out low-quality conversations with no score improvements or more than 6 turns, since we expect that VLMs could learn to quickly and efficiently improve their responses from our data. Finally, we obtain 100K conversations, shown in Fig. 2(a).

### 3.3 FIRE-1M

We use FIRE-100K to fine-tune LLaVA-Next [31] and obtain two models: FIRE100K-LLaVA and FD-LLaVA, which are used to act as the student and the teacher, respectively (training details are shown in Sec. 4). We sample 1M data from $18$ source datasets and generate feedback-refinement conversations via the following steps, as shown in Fig. 3 (b). (1) We feed an image and instruction to the student that generates a response. (2) We feed the image, instruction, ground truth answer, and the response to the teacher that generates feedback. If the score $a$ in the feedback $a \geq 8$ or the number of turns exceeds 3, we stop the conversation; otherwise, go to step (3). (3) We feed the feedback to the
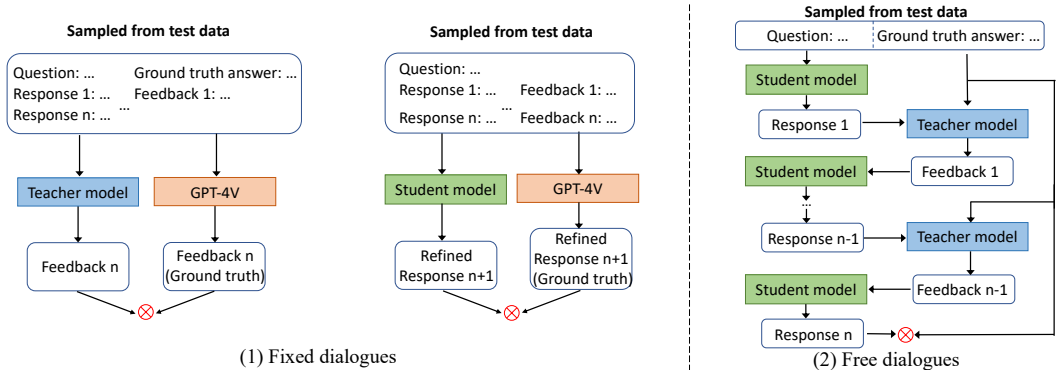
Figure 4: We use two settings to evaluate student and teacher models.

student that generates a refined response and go back to step (2). Finally, we obtain 1M data, shown in Fig. 2(a)

## 3.4 FIRE-Bench

To comprehensively evaluate the feedback-refining ability of VLMs, we introduce FIRE-Bench, containing 11K high-quality feedback-refinement conversations. As shown in Fig. 2(c), FIRE-Bench is derived from 16 source datasets, including 8 seen datasets (test splits) from FIRE-100K and FIRE-1M, as well as 8 new datasets from recently-proposed popular multimodal benchmarks, which is used to evaluate the generalization of the feedback-refining ability across different types of tasks. Similar to FIRE-100K, we sample 11K examples from the data sources and prompt GPT-4V to generate the feedback-refinement conversations.

### 3.4.1 Evaluation Settings

We design two evaluation settings: fixed dialogues and free dialogues to evaluate the performance of the student and teacher models, as shown in Fig. 4.

**Fixed Dialogues.** In fixed dialogues, we evaluate whether the student and teacher models can generate appropriate responses and feedback given the conversation history, and their performance is evaluated by being compared with GPT-4V generated feedback and response, using the BLEU [40] and CIDEr [49] metrics to measure the textual alignment. For the predicted score $\hat{a}_i$ in feedback, we regard the score $a_i$ generated by GPT-4V as the ground truth and adopt *mean absolute error (MAE)*: $MAE = \frac{1}{K} \sum_{k=1}^{K} |a_k - \hat{a}_k|$, where there are $K$ test data totally. The teacher model may fail to follow instructions and does not generate a score in feedback for some cases. Here, we simply set $|a_i - \hat{a}_i| = 10$ for these cases.

**Free Dialogues.** We use a student model and a teacher model to perform free dialogues, and evaluate how fast and how much the student model can improve its answers based on the feedback from the teacher model. The stopping condition for dialogues is that the obtained scores from the teacher model do not increase or exceed a pre-defined threshold (we set 8 in experiments).

We introduce four metrics: average turn (AT), average dialogue refinement (ADR), average turn refinement (ATR), and refinement ratio (RR) for free dialogues.

(1) *Average Turn (AT)*. The AT metric evaluates how fast a VLM could achieve a satisfactory result based on feedback. We measure the number of turns $n_k$ in the conversation to solve the $k$-th data, where VLMs refine their responses until the obtained score exceeds the pre-defined threshold. We set a punishment number as $p = 10$, the maximum number of turns as $n_{max} = 5$. If VLMs fail to obtain a satisfactory score in $n_{max}$ turns, then $n_k = p$. For clearer comparisons with baseline models (*e.g.*, the original LLaVA-Next model), we normalize it according to the AT of the baseline model,

$$AT = \frac{1}{K} \sum_{k=1}^{K} n_k / T_{baseline}, \tag{1}$$

5

where $T_{baseline}$ is the average turn of the baseline model. A smaller value of AT means better performance.

(2) *Average Dialogue Refinement (ADR)*. The ADR metric evaluates how much knowledge VLMs could learn from feedback in a dialogue. In solving the $k$-th data, we use $a_{k,1}$ to denote the obtained score for the initial response, and use $a_{k,n_k}$ to denote the obtained score for the response in the final turn. ADR averages the score improvements of each conversation as

$$ADR = \frac{1}{K} \sum_{k=1}^{K} a_{k,n_k} - a_{k,1}. \tag{2}$$

A larger value of ADR means better performance.

(3) *Average Turn Refinement (ATR)*. ATR evaluates how much knowledge VLMs could learn from feedback in one turn. ATR averages the score improvements in each turn of $K$ samples as

$$ATR = \frac{1}{K} \sum_{k=1}^{K} \frac{1}{n_k - 1} (a_{k,n_k} - a_{k,1}). \tag{3}$$

A larger value of ATR means better performance.

(4) *Refinement Ratio (RR)*. RR measures the proportion of data that have a wrong initial response and a correct final response (*i.e.*, how much data are corrected based on feedback), computed by

$$RR = \frac{1}{K} \sum_{k=1}^{K} \mathbb{1}_{a_{k,n_k} \geq 8} - \mathbb{1}_{a_{k,1} \geq 8}, \tag{4}$$

where $\mathbb{1}_{a_{k,n_k} \geq 8}$ means if $a_{k,n_k} \geq 8$, $\mathbb{1}_{a_{k,n_k} \geq 8} = 1$, and 0 otherwise. A larger value of RR means better performance. Note that, for the $k$-th sample, if $n_k = 1$, we remove it from the K samples to compute AT, ADR, ATR, and RR.

## 3.5 Dataset Analysis

We provide three key statistics: score, turn, and length, for the collected feedback-refinement conversations. **Score.** We show the distribution of initial scores in Fig. 5(a), which reflects the starting state of the conversation. They mainly fall in the interval $[3, 8]$, showing that FIRE covers diverse starting states of conversations. Improved scores per turn are shown in Fig. 5(b), which reflects the learning effect. It ranges from $[2, 8]$, similar to actual situations, where high improvements are obtained in easy cases and small improvements are obtained in hard cases, showcasing the diversity of data. Improved scores per dialogue are shown in Fig. 5(c), and the improvements in most cases are 5-7, demonstrating the data quality of FIRE, where most data have obvious improvements, helping VLMs to efficiently learn to improve their responses. The score distributions of FIRE-100K, FIRE-1M, and FIRE Bench are not completely consistent, making the data more diverse. **Turn.** The turn distribution of conversations is shown in Fig. 5(d). Most conversations have 2-4 turns, indicating an efficient and concise feedback process. This measure suggests that most conversations reach a satisfactory level of refinements. A small number of turns in FIRE informs VLMs to perform effective dialogues. **Length.** The length distributions of responses and feedback are shown in Fig. 5(e) and Fig. 5(f), respectively. Most responses or feedback are less than 100 words. It shows concise dialogues in FIRE, aligning with real-world scenarios where users typically engage in brief exchanges rather than lengthy discussions.

## 4 Model

Our model architecture has the same design as LLaVA-Next-8B [30] that uses CLIP [42] as a frozen image encoder with a two-layer multi-layer perceptron vision-language connector. For the LLM part, we use the same architecture as the LLaMA3-8B [37]. We use LLaVA-Next-8B to initialize the VLMs and use LoRA to fine-tune the LLaVA-Next-8B for a student model and a teacher model.

### 4.1 Student Model

Given an $n$-turn conversation $\{I, q, gt, \{r_i, f_i\}_{i=1}^{n}\}$, we train a student model to fit responses $r_i$ for $i \geq 2$ using the cross-entropy loss,

(a) Score in the first turn     (b) Improved score per turn     (c) Improved score per dialogue

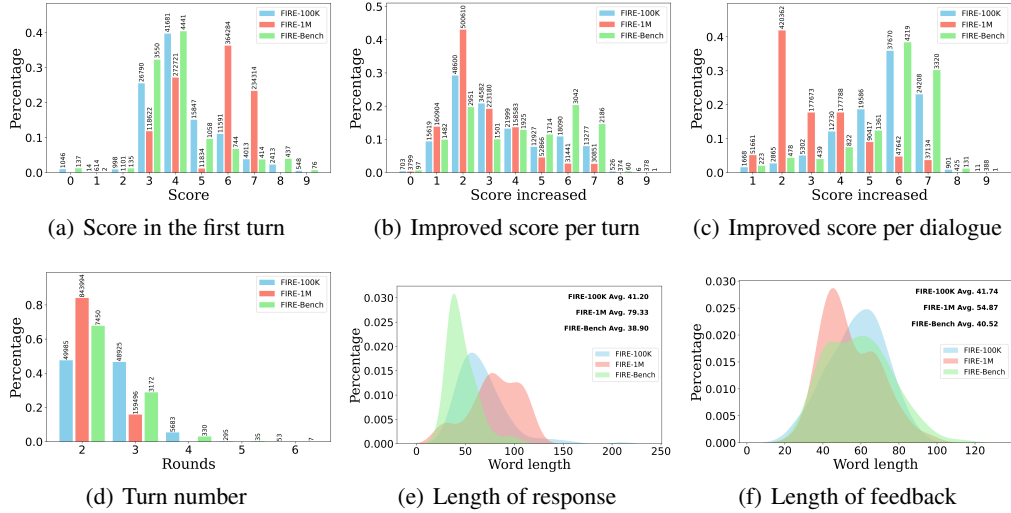(d) Turn number     (e) Length of response     (f) Length of feedback

Figure 5: Data statistics on FIRE-100K, FIRE-1M, FIRE-Bench.

Table 1: Comparisons between LLaVA-Next-8B and FIRE100K-LLaVA on 10 benchmarks. Benchmark names are abbreviated for space limits. GQA [19]; VQAv2 [15];VizWiz [16]; TextVQA [44]; $SQA^I$:ScienceQA-IMG [35]; $LLaVA^W$: LLaVA-Bench-in-the-wild [32];MMB: MMBench [33]; $MME^P$: MME Perception [13]; $MME^C$: MME Cognition [13]; MM-Vet [50].

| Method | GQA | VQAv2 | VizWiz | TextVQA | $SQA^I$ | $LLaVA^W$ | MMB | $MME^P$ | $MME^C$ | MM-Vet |
|---|---|---|---|---|---|---|---|---|---|---|
| LLaVA-Next-8B | **65.9** | 79.0 | 52.0 | **69.8** | **77.3** | **78.5** | 74.4 | 1546 | 331.4 | **44.9** |
| FIRE-LLaVA | 64.5 | **80.9** | **54.3** | 61.0 | 76.8 | 73.4 | **79.3** | **1548** | **340.5** | 38.3 |

$$\min \mathbb{E}_{(I,q,gt,\{r_i,f_i\}_{i=1}^n)\sim\mathbb{D}} \left[ -\sum_{i=2}^n \log P(r_i | I, q, \{r_j, f_j\}_{j=1}^{j=i-1}) \right], \qquad (5)$$

where $\mathbb{D}$ is the used dataset. We first use FIRE-100K as $\mathbb{D}$ to train a student model FIRE100K-LLaVA, then use all training data (FIRE-100K and FIRE-1M) to train a final student model FIRE-LLaVA.

### 4.2 Teacher Model

Given a $n$-turn conversation $\{I, q, gt, \{r_i, f_i\}_{i=1}^n\}$, we train a teacher model to fit the feedback $f_i$ for $i \geq 1$ using the cross-entropy loss,

$$\min \mathbb{E}_{(I,q,gt,\{r_i,f_i\}_{i=1}^n)\sim\mathbb{D}} \left[ -\sum_{i=1}^n \log P(f_i | I, q, gt, \{r_j, f_j\}_{j=1}^{j=i-1}, r_i) \right], \qquad (6)$$

where we use FIRE-100K as $\mathbb{D}$ and obtain the teacher model FD-LLaVA.

## 5 Experiments

We conduct experiments to evaluate both the student and teacher models trained on FIRE. We first provide experimental details and then comprehensively evaluate models in multiple settings.

### 5.1 Experimental Details

**Training Data.** To avoid the catastrophic forgetting issue, we combine the training data in FIRE with the LLaVA-665K [32] (released by Open-LLaVA-1M [6]) to train the student and teacher models.

**Training Details.** In the training process of both the student and teacher models, we freeze the image encoder and the image-language connector, and fine-tune the language decoder using LoRA [18].

Table 2: Results of the student model in fixed dialogues.

| Model | BLEU-1 (↑) | BLEU-2 (↑) | BLEU-3 (↑) | BLEU-4 (↑) | CIDEr (↑) |
|---|---|---|---|---|---|
| LLaVA-Next-8B | 0.33 | 0.23 | 0.17 | 0.13 | 0.60 |
| FIRE-LLaVA | **0.54** | **0.46** | **0.39** | **0.34** | **2.36** |

Table 3: Results of the teacher model in fixed dialogues.

| Model | BLEU-1 (↑) | BLEU-2 (↑) | BLEU-3 (↑) | BLEU-4 (↑) | CIDEr (↑) | MAE (↓) |
|---|---|---|---|---|---|---|
| LLaVA-Next-8B | 0.34 | 0.21 | 0.15 | 0.10 | 0.51 | 1.88 |
| FD-LLaVA | **0.55** | **0.45** | **0.39** | **0.33** | **2.27** | **0.30** |

Table 4: Results in free dialogue over all test data in FIRE.

| Model | AT (↓) | ADR (↑) | ATR (↑) | RR (↑) |
|---|---|---|---|---|
| LLaVA-Next-8B | 1 | 0.97 | 0.41 | 0.25 |
| FIRE100K-LLaVA-8B | 0.92 | 1.27 | 0.55 | 0.34 |
| FIRE-LLaVA-8B | **0.84** | **1.56** | **0.66** | **0.39** |

In the implementation of LoRA, we set the rank as 64 and only apply LoRA on the query and key projection matrices in all attention layers of the language decoder. This setting only involves 0.4% parameters of LLaMA3-8B. We use the AdamW optimizer, where a cosine annealing scheduler is employed, the learning rate is $2e-4$, the batch size is 64, and we train 1 epoch over all data. The training process for a student (or teacher) model requires about 128 A100-80GB GPU hours.

## 5.2 Evaluation in Instruction Following

Considering that fine-tuning VLMs may encounter the catastrophic forgetting problem, we evaluate the instruction-following ability of FIRE-LLaVA, using 10 commonly used multimodal benchmarks, as shown in Tab. 1. Our model achieves comparable performance to the original LLaVA-Next-8B model, showing that we do not compromise the instruction-following ability when learning the feedback-refining ability.

## 5.3 Evaluation in fixed dialogues

We evaluate the performance of FIRE-LLaVA, and FD-LLaVA in fixed dialogues. The evaluation of FIRE-LLaVA is shown in Tab. 2, where we report the results of BLEU-1, BLEU-2, BLEU-3, BLEU-4, and CIDEr. The performance of FD-LLaVA is shown in Tab. 3, where we report the results of BLEU-1, BLEU-2, BLEU-3, BLEU-4, CIDEr, and MAE. We observe that using FIRE, FIRE-LLaVA and FD-LLaVA generates good responses and feedback, having better performance than the original LLaVA-Next-8B model on all metrics. FIRE-LLaVA could well refine the responses, like GPT-4V. FD-LLaVA can generate more accurate feedback, including comments (see BLEU and CIDEr) and scores (see MAE), demonstrating the effectiveness of our teacher model FD-LLaVA that can discover undesirable responses.

## 5.4 Evaluation in the free dialogue

We employ a student model and a teacher model to perform free dialogues. We evaluate LLaVA-Next-8B, FIRE100K-LLaVA, and FIRE-LLaVA as the student model, and use FD-LLaVA to act as the teacher model. We report the average turn (AT), average dialogue refinement (ADR), average turn refinement (ATR), and refinement ratio (RR) on FIRE-Bench. Results are shown in Tab. 4. We observe that a LLaVA model trained on FIRE has improved feedback-refining ability. On the ADR, ATR, and RR metrics, FIRE-LLaVA achieves more than 50% improvements by LLaVA-Next, making an efficient user-agent interaction. Meanwhile, adding FIRE-1M to training data has better performance than only using FIRE-100K, showing the data quality of FIRE-1M.

Table 5: Detailed test results (AT (↓), ADR (↑), ATR (↑), and RR (↑)) on 8 seen source datasets.

| Model | VQAv2 | | | | GQA | | | | TextVQA | | | | ChartQA | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | AT | ADR | ATR | RR | AT | ADR | ATR | RR | AT | ADR | ATR | RR | AT | ADR | ATR | RR |
| LLaVA-Next | 1.00 | 1.45 | 0.42 | 0.40 | 1.00 | 1.51 | 0.51 | 0.43 | 1.00 | 0.91 | 0.34 | 0.26 | 1.00 | 0.71 | 0.39 | 0.25 |
| FIRE100K-LLaVA | 0.86 | 1.83 | 0.55 | 0.54 | 0.81 | 1.93 | 0.63 | 0.58 | 0.95 | 1.20 | 0.49 | 0.33 | 1.07 | 1.03 | **0.56** | 0.27 |
| FIRE-LLaVA | **0.78** | **2.08** | **0.59** | **0.56** | **0.81** | **2.06** | **0.70** | **0.58** | **0.77** | **1.51** | **0.56** | **0.42** | **0.79** | **1.15** | 0.53 | **0.36** |

| Model | DocVQA | | | | DVQA | | | | GEOQA+ | | | | Synthdog | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | AT | ADR | ATR | RR | AT | ADR | ATR | RR | AT | ADR | ATR | RR | AT | ADR | ATR | RR |
| LLaVA-Next | 1.00 | 0.97 | 0.56 | 0.24 | 1.00 | 1.66 | **0.50** | 0.42 | 1.00 | 0.14 | 0.07 | 0.08 | 1.00 | 0.14 | 0.05 | 0.04 |
| FIRE100K-LLaVA | 1.06 | 0.84 | 0.51 | 0.22 | 0.79 | 1.87 | 0.46 | **0.51** | **0.84** | 0.70 | 0.33 | **0.28** | **0.93** | 0.18 | 0.07 | **0.08** |
| FIRE-LLaVA | **0.81** | **1.65** | **0.97** | **0.41** | **0.74** | **1.97** | 0.46 | 0.50 | **0.84** | **0.74** | **0.35** | 0.27 | 0.95 | **0.19** | **0.08** | 0.06 |

Table 6: Detailed test results (AT (↓), ADR (↑), ATR (↑), and RR (↑)) on 8 new source datasets.

| Model | MathVista | | | | MathVerse | | | | MMMU | | | | MME | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | AT | ADR | ATR | RR | AT | ADR | ATR | RR | AT | ADR | ATR | RR | AT | ADR | ATR | RR |
| LLaVA-Next | 1.00 | 0.84 | 0.45 | 0.19 | 1.00 | 0.14 | 0.13 | 0.08 | 1.00 | 0.94 | 0.53 | 0.22 | 1.00 | 1.31 | 0.31 | 0.21 |
| FIRE100K-LLaVA | 0.89 | 1.09 | 0.68 | 0.29 | 0.95 | 0.34 | 0.30 | 0.16 | 0.86 | 1.38 | 0.81 | 0.38 | **0.95** | **2.20** | **0.60** | **0.39** |
| FIRE-LLaVA | **0.83** | **1.36** | **0.77** | **0.34** | **0.93** | **0.65** | **0.49** | **0.17** | **0.80** | **1.73** | **1.05** | **0.41** | 0.96 | 2.04 | 0.57 | 0.36 |

| Model | MM-Vet | | | | SEED-Bench | | | | ScienceQA | | | | LLaVA-wild | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | AT | ADR | ATR | RR | AT | ADR | ATR | RR | AT | ADR | ATR | RR | AT | ADR | ATR | RR |
| LLaVA-Next | 1.00 | 0.80 | 0.31 | 0.13 | 1.00 | 2.30 | 0.56 | 0.48 | 1.00 | 2.81 | 0.70 | 0.56 | 1.00 | 0.45 | 0.19 | 0.03 |
| FIRE100K-LLaVA | 0.97 | 0.99 | 0.48 | 0.23 | 0.83 | 3.18 | 0.75 | 0.68 | 0.98 | 2.95 | 0.78 | 0.62 | 0.99 | 0.79 | 0.33 | **0.12** |
| FIRE-LLaVA | **0.87** | **1.18** | **0.60** | **0.26** | **0.81** | **3.34** | **0.84** | **0.69** | **0.83** | **3.94** | **1.08** | **0.78** | **0.96** | **0.85** | **0.50** | **0.12** |

We also show the detailed results on 8 seen source datasets and 8 new source datasets, as shown in Tab. 5 and Tab. 6, respectively. Our models achieve improvements on both seen and new datasets, showing the generalization of feedback-refining ability across different types of data and tasks.

## 5.5 Visualization

In Fig. 6, we visualize three cases in free dialogues using FIRE-LLaVA and FD-LLaVA. We observe that FIRE-LLaVA can understand the intent behind the feedback and refines its responses according to the feedback on chart understanding, visual concept perception, visual relationship reasoning, and OCR, making efficient user-agent interactions.
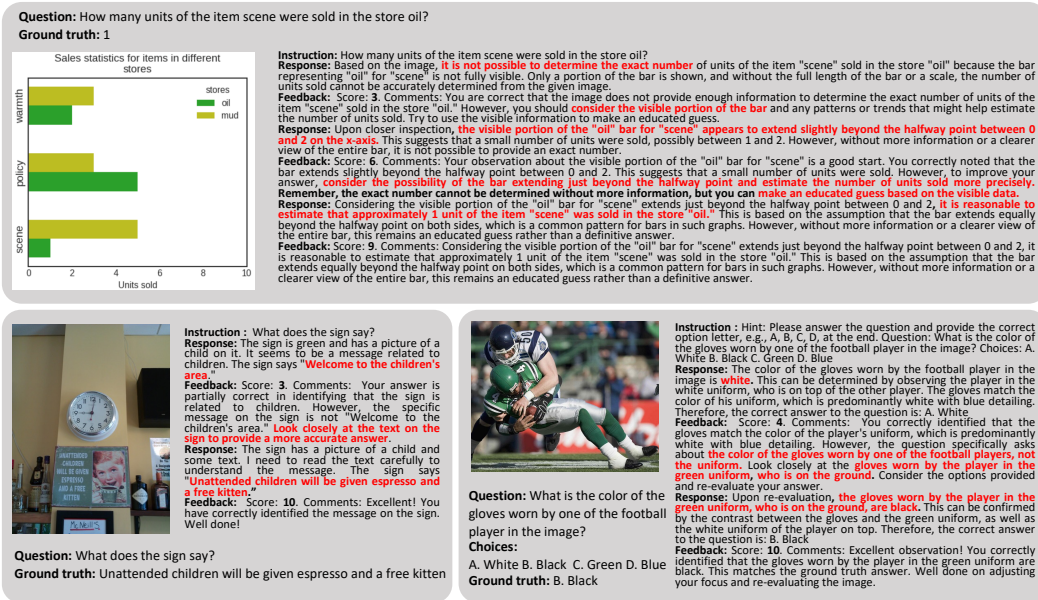


Figure 6: Case study of the feedback-refining ability in our model.

## 6 Conclusion

In this paper, we have presented FIRE, a feedback-refinement dataset with 1.1M multi-turn conversations, which empowers VLMs to refine their responses based on given feedback. Given proper prompts, GPT-4V can produce high-quantity conversations with feedback and responses. Using the 100K GPT-4V generated data as seeds, a student model and a teacher model can freely expand the feedback-refinement data to 1.1M with a similar data quality to GPT-4V. Experiments show that VLMs trained on FIRE have significant improvements in their feedback-refining ability.

**Limitation.** In the current FIRE dataset, the feedback data is limited in the textual form. Practical feedback usually involves diverse multimodal information, such as pointing out image regions. We will further expand FIRE with multimodal feedback data. In addition, although we use a filter process to remove low-quality data, we still cannot completely guarantee the quality of the data. In the future, we will combine human verification with machine verification to improve the quality.

# References

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2

[2] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. In *International Conference on Learning Representations (ICLR)*, 2024. 3

[3] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023. 3

[4] Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. Allava: Harnessing gpt4v-synthesized data for a lite vision-language model. *arXiv preprint arXiv:2402.11684*, 2024. 3

[5] Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023. 3

[6] Lin Chen. Open-llava-next-mix1m dataset, 2024. 7

[7] Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023. 2, 3

[8] Yangyi Chen, Karan Sikka, Michael Cogswell, Heng Ji, and Ajay Divakaran. Dress: Instructing large vision-language models to align and interact with humans via natural language feedback. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3

[9] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Zhong Muyan, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 2, 3

[10] Ching-An Cheng, Andrey Kolobov, Dipendra Misra, Allen Nie, and Adith Swaminathan. Llf-bench: Benchmark for interactive learning from language feedback. *arXiv preprint arXiv:2312.06853*, 2023. 3

[11] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 3

[12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Machine Learning (ICML)*, 2020. 2

[13] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation benchmark for multimodal large language models, 2024. 7

[14] Zhi Gao, Yuntao Du, Xintong Zhang, Xiaojian Ma, Wenjuan Han, Song-Chun Zhu, and Qing Li. Clova: A closed-loop visual assistant with tool usage and update. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3

[15] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6904–6913, 2017. 2, 7

[16] Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3608–3617, 2018. 7

[17] Anwen Hu, Haiyang Xu, Jiabo Ye, Ming Yan, Liang Zhang, Bo Zhang, Chen Li, Ji Zhang, Qin Jin, Fei Huang, et al. mplug-docowl 1.5: Unified structure learning for ocr-free document understanding. *arXiv preprint arXiv:2403.12895*, 2024. 2

10

[18] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations (ICLR)*, 2022. 7

[19] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6700–6709, 2019. 7

[20] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023. 2

[21] Yoonho Lee, Michelle S Lam, Helena Vasconcelos, Michael S Bernstein, and Chelsea Finn. Clarify: Improving model robustness with natural language corrections. *arXiv preprint arXiv:2402.03715*, 2024. 3

[22] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Mimic-it: Multi-modal in-context instruction tuning. *arXiv preprint arXiv:2306.05425*, 2023. 3

[23] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International Conference on Machine Learning (ICML)*, pages 19730–19742, 2023. 2

[24] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning (ICML)*, pages 12888–12900, 2022. 2

[25] Loka Li, Guangyi Chen, Yusheng Su, Zhenhao Chen, Yixuan Zhang, Eric Xing, and Kun Zhang. Confidence matters: Revisiting intrinsic self-correction capabilities of large language models. *arXiv preprint arXiv:2402.12563*, 2024. 3

[26] Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. *arXiv preprint arXiv:2403.18814*, 2024. 3

[27] Jacky Liang, Fei Xia, Wenhao Yu, Andy Zeng, Montserrat Gonzalez Arenas, Maria Attarian, Maria Bauza, Matthew Bennice, Alex Bewley, Adil Dostmohamed, et al. Learning to learn faster from human feedback with language model predictive control. *arXiv preprint arXiv:2402.11450*, 2024. 3

[28] Yuan-Hong Liao, Rafid Mahmood, Sanja Fidler, and David Acuna. Can feedback enhance semantic grounding in large vision-language models? *arXiv preprint arXiv:2404.06510*, 2024. 3

[29] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Mitigating hallucination in large multi-modal models via robust instruction tuning. In *International Conference on Learning Representations (ICLR)*, 2024. 3

[30] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2024. 2, 6

[31] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 2, 4

[32] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 2, 3, 7

[33] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*, 2023. 7

[34] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *International Conference on Learning Representations (ICLR)*, 2024. 2

[35] Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pages 2507–2521, 2022. 7

[36] Ahmed Masry, Xuan Long Do, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 2263–2279, 2022. 2

[37] Meta. Llama3, 2024. Accessed: 2024-06-01. 6

[38] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual question answering by reading text in images. In *International conference on document analysis and recognition (ICDAR)*, pages 947–952, 2019. 2

[39] OpenAI. Gpt-4v(ision) system card. 2023. 2

[40] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 311–318, 2002. 5

[41] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. In *International Conference on Learning Representations (ICLR)*, 2024. 3

[42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, pages 8748–8763, 2021. 2, 6

[43] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8317–8326, 2019. 2

[44] Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8317–8326, 2019. 7

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

[46] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 2

[47] Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Haitao Mi, and Dong Yu. Toward self-improvement of llms via imagination, searching, and criticizing. *arXiv preprint arXiv:2404.12253*, 2024. 3

[48] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *ArXiv*, abs/2307.09288, 2023. 2

[49] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4566–4575, 2015. 5

[50] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. In *International Conference on Machine Learning (ICML)*, 2024. 7

[51] Pan Zhang, Xiaoyi Dong Bin Wang, Yuhang Cao, Chao Xu, Linke Ouyang, Zhiyuan Zhao, Shuangrui Ding, Songyang Zhang, Haodong Duan, Hang Yan, et al. Internlm-xcomposer: A vision-language large model for advanced text-image comprehension and composition. *arXiv preprint arXiv:2309.15112*, 2023. 3

[52] Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. Llavar: Enhanced visual instruction tuning for text-rich image understanding. *arXiv preprint arXiv:2306.17107*, 2023. 3

[53] Bo Zhao, Boya Wu, and Tiejun Huang. Svit: Scaling up visual instruction tuning. *arXiv preprint arXiv:2307.04087*, 2023. 3

[54] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. In *International Conference on Learning Representations (ICLR)*, 2024. 2, 3

## Checklist

1. For all authors...

   (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes] The contributions and scope are claimed in the abstract and introduction.

   (b) Did you describe the limitations of your work? [Yes] Yes, we have discussed the limitations in Sec. 6.

   (c) Did you discuss any potential negative societal impacts of your work? [Yes] Yes, we have discussed the potential negative societal impacts of our work in the supplementary materials.

   (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes] We ensure that our paper is fully conforming to the ethics review guidelines.

2. If you are including theoretical results...

   (a) Did you state the full set of assumptions of all theoretical results? [N/A] This paper does not include theoretical results.

   (b) Did you include complete proofs of all theoretical results? [N/A] This paper does not include theoretical results.

3. If you ran experiments (e.g. for benchmarks)...

   (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] The code and data can be found in the URL, `mm-fire.github.io`. The detailed setting of experiments can be found in Sec. 5. The prompts for generating the conversations are shown in the supplementary materials.

   (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] All the training details can be found in Sec. 5.

   (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] The error bars for the student model and the teacher model are discussed in the supplementary materials.

   (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] The computing resources are included in Sec. 5.

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

   (a) If your work uses existing assets, did you cite the creators? [Yes] We cited the creators of all the existing assets used in this paper.

   (b) Did you mention the license of the assets? [Yes] We mentioned the license of the assets in the supplementary materials.

   (c) Did you include any new assets either in the supplemental material or as a URL? [No] We did not include new assets in the supplemental material or as a URL.

   (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A] All the data sources of this paper are public datasets.

   (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A] All the data sources of this paper are public datasets, which did not include any personally identifiable information.

5. If you used crowdsourcing or conducted research with human subjects...

   (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] No crowdsourcing and human subjects are involved in this paper.

(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] No crowdsourcing and human subjects are involved in this paper.

(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] No crowdsourcing and human subjects are involved in this paper.