# InterFeedback: Unveiling Interactive Intelligence of Large Multimodal Models with Human Feedback

Henry Hengyuan Zhao\*, Wenqi Pei\*, Yifei Tao\*, Haiyang Mei, Mike Zheng Shou<sup>†</sup>

Show Lab, National University of Singapore

#### **Abstract**

Existing benchmarks do not test Large Multimodal Models (LMMs) on their interactive intelligence with human users, which is vital for developing general-purpose AI assistants. We design InterFeedback, an interactive framework, which can be applied to any LMM and dataset to assess this ability autonomously. On top of this, we introduce InterFeedback-Bench that evaluates interactive intelligence using two representative datasets, MMMU-Pro and MathVerse, to test 10 different open-source LMMs. Additionally, we present InterFeedback-Human, a newly collected dataset of 120 cases designed for manually testing interactive performance in proprietary models such as OpenAI-o1 and Claude-3.5-Sonnet. Our evaluation results indicate that even the state-of-the-art LMM, OpenAIo1, struggles to refine its responses based on human feedback, achieving an average score of less than 50%. Our findings point to the need for methods that can enhance LMMs' capabilities to interpret and benefit from feedback.

#### 1 Introduction

In this paper, we are curious about the question "Can Large Multimodal Models evolve through Interactive Human Feedback?" It is central to developing general-purpose AI assistants with Large Multimodal Models (LMMs). While these models show exceptional performance on tackling multimodal tasks directly, their ability to interact with humans remains largely unknown. We argue that an LMM functioning as the general assistant should possess two capabilities: 1) exceptional problemsolving ability and 2) the ability to improve itself through feedback (e.g., human feedback, execution results). In this work, we focus on the latter capability, which has been rarely examined in existing benchmarks.

Humans are highly adaptive, continuously refining their skills through feedback-a fundamental process for acquiring knowledge and solving problems. Likewise, advanced LMM models should be designed to learn from feedback, ensuring better alignment with real-world needs and enhancing their problem-solving capabilities in Human-AI Interaction (HAI). Recently, a surge of large multimodal models (LMMs) (OpenAI, 2023; Wang et al., 2024; Deitke et al., 2024; Zhao et al., 2024b; Li et al., 2024a; Zhao et al., 2024a; Chen et al., 2024b) have been developed to handle various tasks, including general vision-language understanding (Liu et al., 2023b; Li et al., 2023), expert-level multimodal understanding (Yue et al., 2024a,b), and scientific reasoning (Lu et al., 2022, 2024; Zhang et al., 2024). However, these LMMs are primarily tested in a static way, overlooking their great potential in an interactive process such as interactive coding (Jimenez et al., 2024; Yang et al., 2025), computer usage (Zhao et al., 2025; Lin et al., 2024; Gao et al., 2024; Xie et al., 2024), and clinical reasoning (Li et al., 2024e). Consequently, the interactive intelligence of LMMs remains largely unexplored, and the development of a standard benchmark for evaluating their interactive intelligence remains an open challenge.

The key challenge in evaluating the interactive intelligence of LMMs is the automatic model tests. In practice, for the same query, different LMMs often produce varied responses, necessitating that humans offer tailored feedback for each conversation round. To address this issue, we propose InterFeedback a straightforward problem-solving framework that enables any LMM to tackle multimodal tasks interactively by leveraging proprietary models such as GPT-40 (OpenAI, 2023) to simulate humans, inspired in previous studies (Yao et al., 2025; Chen et al., 2024a; Yoon et al., 2024; Luo et al., 2024).

On top of this framework, we present

<sup>\*</sup>Equal Contribution.

<sup>&</sup>lt;sup>†</sup>Corresponding author.

**InterFeedback-Bench**, a benchmark designed to comprehensively evaluate LMMs for two purposes: 1) the ability to interactively solve problems and 2) the capability of interpreting the feedback to improve themselves. We demonstrate with two challenging pre-existing datasets: MMMU-Pro (Yue et al., 2024b) and Mathverse (Zhang et al., 2024). Additionally, for a more in-depth investigation, we conduct human evaluation on four closed-source proprietary models: GPT-4o (OpenAI, 2023), OpenAI-o1 (OpenAI, 2024), Claude-3.5-Sonnet (Anthropic, 2024), and Gemini-2.0 (Gemini, 2025) with a trained user acting as the feedback provider. Finally, we manually collected a dataset InterFeedback-Human containing 120 samples for this assessment.

Our experimental results reveal several compelling insights: 1) Interactive process could improve the performance of most LMMs in solving challenging problems; 2) Existing LMMs exhibit suboptimal performance in interpreting and incorporating feedback; 3) Accuracy result may not truly reflect the model's capability to improve itself from feedback; 4) High-quality feedback is essential, as subpar feedback can degrade performance even more than a simple binary (0/1) correctness signal; 5) LMM may not truly reasoning, we find out that LMMs resort to guessing answer even on a simple question according to human. These findings point to the need for methods that can enhance the LMM's capability to interpret and benefit from feedback. In summary, our contributions are:

- We take the first step toward exploring the interactive intelligence of LMMs in improving themselves through human feedback.
- We propose a straightforward and extensible framework InterFeedback which allows any LMM to interactively solve problems.
- We construct InterFeedback-Bench, a novel and universal benchmark for assessing the ability of interactive problem-solving of LMMs.
- We conduct comprehensive evaluations and indepth analysis, providing several compelling insights for future model alignment.

#### 2 Related Work

Large Multimodal Models. The LLaVA-series works (Liu et al., 2023a, 2024a,b; Li et al.,

2024a) demonstrate that training with supervised fine-tuning (SFT) multimodal data and expand the vision lens would produce compatible multimodal reasoning ability. By adopting a largescale image-text corpus for instruction tuning, Qwen2-VL (Wang et al., 2024), CogVLM (Wang et al., 2023), InternVL2 (OpenGVLab, 2024) have achieved exceptional performance on various multimodal abilities. Moreover, Molmo (Deitke et al., 2024) proposes to train an LMM from scratch with only the human-annotated data. Unlike these large models, MiniCPM-V (Yao et al., 2024) and Phi-3.5-Vision (Abdin et al., 2024) propose to train lightweight yet SOTA LMMs. Despite their exceptional performance on multimodal benchmarks of varying difficulty, such as MMMU-Pro (Yue et al., 2024b) and MathVista (Lu et al., 2024), it remains unclear how well these LMMs demonstrate interactive intelligence in Human-AI Interaction scenarios. In this paper, we conduct the evaluation of these LMMs to explore this basic yet vital capability (i.e., evolve through interactive human feedback).

Multimodal Benchmarks. Traditional visionlanguage benchmarks focus on visual question answering (Goyal et al., 2017), image captioning (Chen et al., 2015), as well as other benchmarks for specialized scenarios such as scene text understanding (Singh et al., 2019), commonsense reasoning (Zellers et al., 2019), outside knowledge (Marino et al., 2019; Schwenk et al., 2022). The recent development of LMM posts a strong need for modernized multimodal benchmarks (Zhao et al., 2025; Liu et al., 2023b; Li et al., 2023; Yu et al., 2023; Yue et al., 2024a; Lu et al., 2024; Zhang et al., 2024; Li et al., 2024d) such as MMBench (Liu et al., 2023b), MMMU-pro (Yue et al., 2024b), and MathVerse (Zhang et al., 2024) which involve comprehensively evaluating current LMMs on various multimodal abilities. However, these benchmarks primarily focus on static testing processes, overlooking the interactive testing process that is vital in human-AI interaction scenarios.

Human-AI Interaction. Investigating how humans and AI systems communicate and collaborate is critical for shaping applications such as virtual assistants (Virvou, 2022), personalized recommendations (Dodeja et al., 2024), autonomous vehicles (Zhang et al., 2021), and healthcare diagnostics (McKinney et al., 2020). Recent LLMs-driven techniques, such as memory (Park et al., 2023) and iterative (Zhang et al., 2023) mechanisms offer expertlevel collaboration. While LMMs (Deitke et al.,

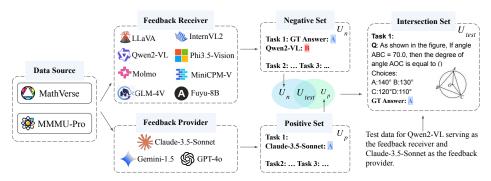


Figure 1: Overview of the test data construction process for InterFeedback-Bench. For each LMM serving as the feedback receiver, we process each instance from a target dataset (e.g., MathVerse) and collect the error cases to form a negative set. The feedback provider then processes the same instances to build a positive set. Finally, we curate test data by selecting the intersection of both sets.

2024; Wang et al., 2024) excel in multimodal tasks, their potential for HAI problem-solving (Yang et al., 2025; Li et al., 2024e) remains underexplored. By offering a unified framework and meticulously curated data, our InterFeedback-Bench enables evaluation of LMMs on these capabilities and lays a foundation for advancing multimodal HAI problem-solving.

User Stimulation with LLM. Recently, previous work in order to build multi-agent system (Khan et al., 2024), stimulate human-AI interaction (Yao et al., 2025), evaluate LMMs in video analysis (Luo et al., 2024), stimulate real users in a web shopping scenario (Chen et al., 2024a), evaluate the conversational recommender systems (Yoon et al., 2024) determine whether to use LLM or LMM to stimulate the user. However, previous works have overlooked the importance of ensuring the reliability of LLMs or LMMs that are used to stimulate the users. In this paper, we curate test data by selecting only the samples that LMMs correctly address, minimizing unreliable interaction results.

#### 3 InterFeedback-Bench

In this section, we first introduce the automated interactive benchmarking process in Section 3.1. We begin by formulating the concept of interactive problem-solving, followed by a discussion of the data curation process. We then present the proposed interactive framework, InterFeedback. Next, in Section 3.2, we elaborate on the human benchmarking component, detailing the data collection and the proposed feedback providing strategy.

## 3.1 Automated Interactive Benchmarking

#### 3.1.1 Formulation

The InterFeedback-Bench formalizes the interactive problem-solving process with feedback in a partially observable Markov decision process (POMDP)  $(S, \mathcal{O}, \mathcal{A}, \mathcal{T}, \mathcal{R})$  with state space S, observation  $\mathcal{O}$ , action space  $\mathcal{A}$ , transition function  $\mathcal{T}$ :  $\mathcal{S} \times \mathcal{A} \to \mathcal{S}$ , and reward function  $\mathcal{R}: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ . In our setting, given a natural language question q(e.g., Please select the sitting camel that is being led and facing right) and the input image v, the model first gets the observation  $o_t \in \mathcal{O}$  from the state  $s_t \in \mathcal{S}$  in the execution environment and then generate the action  $a_t \in \mathcal{A}$ . The  $a_t$  is the response from models in natural language. The reward function  $\mathcal{R}$ :  $\mathcal{S} \times \mathcal{A} \rightarrow \{0,1\}$  here returns a binary value indicating the task correctness status. It is implemented by the exact match: returning 1 if the predicted answer exactly matches the ground-truth, and 0 otherwise. The observation  $o_t$  includes both the correctness signal from the reward function and the feedback from the humans.

#### 3.1.2 Data Curation

**Data sources.** To ensure the quality and difficulty of multimodal tasks, inspired by previous benchmarks demonstrated on pre-existing datasets (Yang et al., 2023; Li et al., 2024c), we choose to test LMMs on two challenging datasets: MathVerse (Zhang et al., 2024) and MMMU-Pro (Yue et al., 2024b). MathVerse is a visual math benchmark that includes various mathematical problems, and 3,940 samples (testmini set) are used in our work. MMMU-Pro is a comprehensive multimodal benchmark and we use 1,730 expert-level questions (single image mode). Both datasets are challenging even for the model GPT-40 which achieves only

64.7% accuracy on MMMU-Pro (Standard 4 options).

**Data selection process.** We choose to use proprietary LMMs, such as GPT-40, for stimulating the humans to give feedback mimicking human-AI interactions. The primary challenge, however, is ensuring that the feedback generated by these models is reliable, as even models like GPT-40 and Claude-3.5-Sonnet still do not perform correctly on all test samples. Therefore, we construct the test data by selecting the intersection set that feedback provider  $M_p$  solves correctly while  $M_r$  does not, as shown in Figure 1. Specifically, the pipeline includes three parts: 1) feedback receiver LMM locally-running; 2) feedback provider LMM APIcalling; and 3) intersection set selection. Such a data construction process leads to each tested LMM having a different test data set.

Specifically, given a test dataset D, we begin by having the feedback receiver model  $M_r$  process every instance in D to produce a negative set  $U_n$  consisting of tasks it fails to solve correctly. Next, the feedback provider model  $M_p$  processes the same dataset to generate a positive set  $U_p$  comprising tasks it solves correctly. We then define  $U_{\text{test}}$  as the intersection of  $U_n$  and  $U_p$ , i.e.,

$$U_{\text{test}} = U_n \cap U_p,$$

which means that  $U_{\rm test}$  contains tasks that  $M_p$  solves correctly but  $M_r$  does not. This approach ensures that the feedback generated by  $M_p$  is both relevant and reliable.

#### 3.1.3 InterFeedback Framework

To enable an interactive problem-solving process, we propose a new straightforward framework Inter-Feedback. It includes two roles: feedback receiver  $M_r$  and feedback provider  $M_p$ , as shown in Figure 2. The feedback receiver is the candidate LMMs (e.g., Qwen2-VL) ready for the benchmark and the feedback provider is the SOTA LMM (e.g., GPT-40) for providing the pertinent feedback in each time step in place of a human. Consider at timestep t, the output of  $M_r$  is  $a_t$ , and the feedback provider  $M_p$  has to follow the policy that provides the feedback  $f_t$  from the mapping :  $F(a_t, s_t) \rightarrow f_t$ . The  $s_t$  denotes the correctness signal from the verification process via the reward function. We record the model outputs for the final evaluation.

**Feedback types.** Additionally, we propose a simplified feedback mechanism that only indicates correctness (i.e., correct or incorrect), without a

detailed explanation. In summary, we evaluate the models using two feedback types: *Detail* and *Simple*. The *Detail* feedback comprises both *Simple* feedback and detailed LMM-generated explanation.

#### 3.2 Human Benchmarking

In the previous section, we employed proprietary LMMs as feedback providers. Naturally, how well do these models perform when receiving feedback? We begin to assess the proprietary LMMs with a human-in-the-loop process. The feedback provider  $M_p$  is a trained user who fully understands all the questions in the newly curated dataset InterFeedback-Human. The feedback receiver  $M_r$  is the proprietary LMMs including OpenAI-o1, GPT-40, Gemini-2.0, and Claude-3.5-Sonnet. This evaluation aims to assess how effectively these proprietary models can serve as assistants in a human-AI interaction system.

#### 3.2.1 Data Collection.

We gather challenging data examples across diverse domains: visual logic, mathematics, and coding. These were selected to probe the cognitive depth of the models, especially when confronted with complex, multi-step reasoning problems. The visual logic data is manually collected from publicly available resources. The emphasis on visual logic tasks reflects the growing demand for models to handle image-based reasoning challenges, such as pattern recognition (Wei et al., 2025) (e.g., determining the next shape in a sequence) and character-based logic (e.g., interpreting transformations between symbols). We also collect the multimodal mathematics data from the existing dataset MathVerse (Zhang et al., 2024) and the multimodal expertlevel data from MMMU-Pro (Yue et al., 2024b). Additionally, we also involve the natural language task into InterFeedback-Human to analyze such capability in the NLP area.

In summary, InterFeedback-Human encompasses a total of 120 tasks distributed across the five task types: 80 visual logic tasks, 10 mathematical logic tasks (sampled from NuminaMath (Li et al., 2024b)), 10 coding tasks (sampled from CodeComprehension (Imbue, 2024)), 10 MMMU-Pro tasks, and 10 MathVerse tasks.

## 3.2.2 Hierarchical Feedback

We design a hierarchical feedback generation scheme to gradually increase the information in-

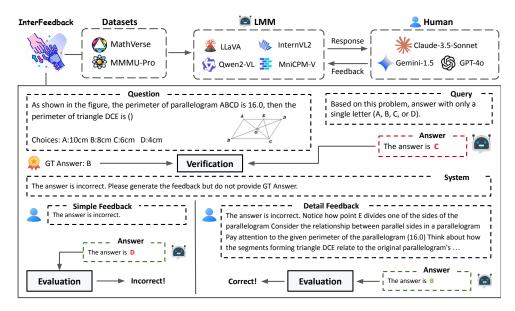


Figure 2: Overview of the proposed framework InterFeedback for assessing an LMM's ability to improve itself through feedback. The model interacts with humans to progressively solve a problem, and after each conversation round, we verify the correctness of the answer. If the answer is incorrect, an LMM-stimulated human will provide constructive feedback. We implement two types of feedback to investigate the behavior of LMMs.

tensity. Specifically, we ask the human to give the following three-level feedback:

- Level 1: Provide a basic and simple description that leads to the correct answer.
- Level 2: Provide an expanded explanation that leads to the correct answer.
- Level 3: The correct answer is <u>GT Answer</u>.
   Provide a comprehensive and detailed explanation that leads to the correct answer.

Since most of our questions have four options, giving more than three rounds of feedback might let the model guess the answer by elimination rather than by reasoning. For example, if the correct answer is A and the model has already given B, C, and D, a third round of feedback is unnecessary. Therefore, we directly provide the <u>GT Answer</u> in Level 3 feedback prompts to test the models' ability to explain their thinking process.

#### 3.2.3 Evaluation Integration

To ensure fairness and consistency in our evaluation, we engaged only one experienced user. Since human-in-the-loop feedback is inherently subjective, involving multiple participants could introduce variability due to differences in background and expertise. This approach helps maintain the reliability of the relative performance comparisons across candidate LMMs.

# 4 Experiments

#### 4.1 Experiment Setup

**Evaluation Models.** We evaluate the performance of foundation models served as the feedback receiver  $M_r$  across 12 representative LLaVA-1.5-7B (Liu et al., 2024a), LMMs: LLaVA-1.6-7B (Liu et al., 2024b) (Mistral-7B), LLaVa-OneVision-7B (Li et al., 2024a) (Qwen2-7B (Yang et al., 2024)), Owen2-VL-7B (Wang et al., 2024), Qwen2.5-VL-7B (Team, 2025b), GLM-4V-9B (Wang et al., 2023), InternVL2 (OpenGVLab, 2024), Molmo (Deitke et al., 2024), MiniCPM-V (Yao et al., 2024), Phi-3.5-Vision (Abdin et al., 2024), Fuyu-8B (Bavishi et al., 2023), and Seed-1.5-VL-Thinking <sup>1</sup> (Team, 2025a). The feedback provider  $M_p$  includes the three best available models from three model (gpt-4o-2024-08-06), families: OpenAI Gemini (Gemini-1.5-Pro), and Claude (claude-sonnet-4-20250514).

**Evaluation Metrics.** In addition to the **Accuracy** metric, we leverage the **Correction Rate**, defined as the percentage of corrected answers of all erroneous samples. Let N denote the total number of samples,  $N_e$  the number of erroneous samples, and  $N_c$  the number of samples that have been corrected. The Accuracy and Correction Rate metrics

<sup>&</sup>lt;sup>1</sup>doubao-1-5-thinking-vision-pro-250428

Model				GPT-4c	)	Gemini-1.5-Flash				Claude-Sonnet-4	
model	Acc (%)	# Neg	# Test	Detail (%)	Simple (%)	# Test	Detail (%)	Simple (%)	# Test	Detail (%)	Simple (%)
				No	on-Thinking N	Models					
LLaVa-OneVision-7B	25.6	2933	373	36.2	18.0	428	29.0	15.7	820	38.3	23.8
Molmo-7B	25.6	2931	452	55.1	52.0	507	36.5	38.9	987	15.3	34.3
MiniCPM-V	16.2	3301	552	28.4	20.3	741	16.6	25.4	1195	5.3	12.1
GLM-4V-9B	20.2	3146	440	38.6	28.2	568	30.1	29.9	1015	22.9	21.9
Phi3.5-Vision-4.2B	19.0	3192	534	36.1	33.7	579	31.3	33.7	1045	21.1	24.2
LLaVa-1.5-7B	13.5	3409	763	23.2	14.3	678	18.0	14.7	1256	3.3	5.8
LLaVa-1.6-Mistral-7B	14.8	3357	549	41.0	35.9	661	5.9	5.9	1212	14.8	17.7
Fuyu-8B	21.8	3083	582	24.1	19.8	635	15.0	12.9	1187	17.9	15.5
InternVL2-8B	38.1	2440	379	49.6	41.2	375	48.8	44.4	547	21.4	26.7
Qwen2-VL-7B	22.5	3052	295	66.8	72.2	470	41.9	44.9	774	34.4	35.8
Qwen2.5-VL-7B	31.5	2698	266	69.2	62.4	350	45.4	42.6	1521	46.8	43.9
Thinking Models											
Seed-1.5-VL-Thinking	47.4	2072	73	67.1	63.0	70	64.3	58.6	474	88.6	90.5

Table 1: Correction Rate Results of three Feedback Providers on MathVerse Dataset. Acc (%): The average accuracy of MathVerse's *testmini* set. (Calculated by our prompt template.) The results are tested by ourselves. # Neg: The number of negative samples produced by the model. # Test: The total number of test samples evaluated. Detail (%): correction rate of using LMM-generated feedback. Simple (%): correction rate of using simple feedback (0 or 1).

Model				GPT-4c	)	Gemini-1.5-Flash Claude-Son			net-4		
110401	Acc (%)	# Neg	# Test	Detail (%)	Simple (%)	# Test	Detail (%)	Simple (%)	# Test	Detail (%)	Simple (%)
				No	on-Thinking N	Models					
LLaVa-OneVision-7B	47.1	915	312	31.7	15.7	333	35.4	18.6	539	42.2	30.6
Molmo-7B	43.8	973	362	51.7	48.9	383	41.5	43.1	593	19.7	33.9
MiniCPM-V	38.1	1071	410	27.3	23.7	503	21.5	21.7	688	7.0	15.3
GLM-4V-9B	46.0	935	327	38.8	30.0	359	38.7	31.5	577	27.6	23.6
Phi3.5-Vision-4.2B	43.2	983	366	44.3	42.3	396	40.9	39.6	611	31.8	31.1
LLaVa-1.5-7B	36.5	1099	506	31.9	12.3	470	20.0	16.0	720	8.6	11.8
LLaVa-1.6-Mistral-7B	38.8	1058	432	46.1	36.1	429	14.7	14.7	682	27.3	25.5
Fuyu-8B	34.1	1140	481	6.0	8.7	1140	3.7	3.5	768	10.2	8.7
InternVL2-8B	45.7	939	343	50.1	41.4	329	57.1	50.2	435	23.7	32.9
Qwen2-VL-7B	48.1	898	268	50.4	44.8	322	39.4	37.6	525	35.6	33.7
Qwen2.5-VL-7B	50.0	865	839	39.6	36.1	323	44.9	39.0	839	39.6	36.1
Thinking Models											
Seed-1.5-VL-Thinking	94.2	101	20	80.0	70.0	31	64.5	64.5	50	74.0	70.0

Table 2: Correction Rate Results of three Feedback Providers on MMMU-Pro Dataset. We test models on a single image setting of MMMU-Pro. The accuracy and number of correct answers for the three feedback providers are GPT-40 with 42.6% (737), Gemini-1.5-Flash with 47.1% (815), and Claude-3.5-Sonnet with 52.3% (905).

can be formulated as follows:

$$Accuracy = \frac{N - N_e}{N} \times 100\%, \qquad (1)$$

Correction Rate = 
$$\frac{N_c}{N_e} \times 100\%$$
. (2)

Implementation Details. We set the temperature to 0 for all tested models and API models. The image resolution of the Qwen2-VL model we restrict to  $512 \times 512$  to avoid the memory exceeded error. All evaluations were conducted on two NVIDIA RTX A6000 GPUs. To ensure the reliability of results, we obtain the intersection set for both the feedback receiver and provider models that are able to output the correct answer format. Based on our preliminary experiments, we limited the interactive benchmarking to a single round. This decision is driven by two observations: most

models fail to provide correct answers in subsequent rounds, and multiple rounds tend to lead to answer guessing, which undermines the reliability of quantitative evaluation.

**Feedback Types.** As introduced in Section 3.1, we employ proprietary LMMs to stimulate the human to provide pertinent feedback at each conversation round. Additionally, we propose a simplified feedback mechanism that only indicates correctness (i.e., correct or incorrect), without a detailed explanation. In summary, we evaluate the models using two feedback types: *Detail* and *Simple*. The *Detail* feedback comprises both *Simple* feedback and detailed LMM-generated feedback.

# **4.2** Experimental Analysis on Interactive Benchmarking

To thoroughly investigate the ability of LMMs to integrate feedback and improve their problem-solving performance, we present evaluation results for various models on two datasets—MathVerse (Zhang et al., 2024) in Table 1 and MMMU-Pro (Yue et al., 2024b) in Table 2, respectively. Below, we provide a detailed discussion of key findings.

Whether interactive process improves the performance of LMMs? Yes. As demonstrated in both tables, integrating our proposed framework InterFeedback enables most models to benefit from feedback provided by SOTA LMMs, such as GPT-40 and Claude-Sonnet-4. Notably, even the weaker model Fuyu-8B sees 24.1% of its erroneous samples corrected through GPT-40's feedback.

**Current LMMs struggle to enhance perfor**mance through feedback. As shown in the tables, most LMMs are unable to correct all erroneous samples, even when provided with feedback from state-of-the-art proprietary models such as Claude-Sonnet-4 and GPT-4o. For example, consider the two SOTA open-source models, Qwen2-VL-7B and Molmo. Qwen2-VL-7B achieves a 66.8% correction rate on the MathVerse dataset with GPT-4o's feedback, but only a 50.4% correction rate on the MMMU-Pro dataset. Similarly, Molmo-7B attains correction rates of 55.1% and 51.7% on the MathVerse and MMMU-Pro datasets, respectively. Overall, the correction rates for the rest of the models remain below 50%. This suggests that even with constructive feedback from advanced LMMs, current models struggle to enhance performance through feedback generally.

Accuracy result may not truly reflect the model's capability to improve itself from feedback. As shown in Table 1, although InternVL2-8B achieves a higher accuracy (38.1%), its correction rate is only 49.6%. In contrast, Qwen2-VL-7B, with a lower accuracy of 22.5%, attains the highest correction rate of 66.8% when using GPT-4o's feedback. Similarly, Molmo-7B surpasses InternVL2-8B in correction rate despite having lower accuracy. On the MMMU-Pro dataset (see Table 2), LLaVA-OneVision-7B records 47.1% but only a 31.7% correction rate, which is lower than that of several models that have inferior accuracy (e.g., InternVL2-8B, Molmo-7B, GLM-4v-9B, and Phi3.5-Vision-4.2B). This inconsistency between initial answering ability and self-improvement capability indicates that evaluating models solely on accuracy may not fully capture their true potential.

Simple feedback also enhances performance. In addition to using detailed LMM-generated feedback, we evaluated models with binary (0/1) feedback that simply indicates the correctness of their current response. Surprisingly, the results show that all models benefit from this simple feedback mechanism. This suggests that while LMMs have the inherent potential to generate correct answers, they may require additional prompting techniques to fully harness their problem-solving capabilities.

LMM-generated feedback is not always better than simple feedback. By comparing the results obtained using *Detail* feedback from GPT-40 with those using *Simple* binary feedback, we observe that most models perform better with detailed feedback. For example, on the MathVerse dataset, LLaVA-OneVision-7B achieves 36.2% with detailed feedback versus 18.0% with binary feedback; InternVL2-8B increases from 41.2% to 49.6%; and MiniCPM-V increases from 20.3% to 28.4%. However, Qwen2-VL scores 66.8% with detailed feedback and 72.2% with simple feedback. Similarly, on the MMMU-Pro dataset, Fuyu-8B performs worse with detailed feedback (6.0% vs. 8.7%).

The quality of feedback is crucial: low-quality feedback can degrade performance more than simply providing binary (0/1) feedback. We compare the feedback provided by GPT-40 and Gemini-1.5-Flash on the challenging MathVerse dataset, where most models achieve accuracies below 30%, highlighting the difficulty of its problem instances. We find that delivering simple binary feedback that merely indicates the correctness of the tested model's output can outperform LMM-generated detailed feedback (Gemini-1.5-Flash). Specifically, the correction rates using simple feedback exceed those with detailed feedback for several models: Molmo-7B (38.9% vs. 36.5%), MiniCPM-V (25.6% vs. 16.6%), Phi3.5-Vision-4.2B (33.7% vs. 31.3%), and Qwen2-VL-7B (44.9% vs. 41.9%).

# 4.3 Experimental Analysis on Human Benchmarking

In this section, we will introduce the human evaluation results of several well-known closed-source families: OpenAI (GPT-40, OpenAI-o1), Claude (Claude-3.5-Sonnet-20241022), and Gemini (Gemini-2.0-Flash-Exp).

Overall Accuracy Results. In Table 3: (1)

Model	Visual Logic	MMMU-Pro	MathVerse	Math <sup>Text</sup>	Coding <sup>Text</sup>	Average
Gemini-2.0	21.3	50.0	70.0	50.0	50.0	32.5
Claude-3.5	37.5	60.0	80.0	70.0	70.0	48.3
OpenAI-o1	28.8	60.0	90.0	90.0	90.0	46.7
GPT-40	25.0	70.0	80.0	60.0	50.0	38.3

Table 3: **Human Evaluation Results across LMMs on InterFeedback-Human.** Math<sup>Text</sup> and Coding<sup>Text</sup> represent two text-only task categories. The scores represent the average percentage of correct samples among all samples.

Model	# Round	Visual Logic	MMMU-Pro	MathVerse	Math <sup>Text</sup>	Coding <sup>Text</sup>	Average
	1	38.1	20.0	33.3	0.0	80.0	37.0
Gemini-2.0	2	20.6	0.0	33.3	20.0	20.0	19.8
	3	41.3	80.0	33.3	80.0	0.0	43.2
	1	38.0	0.0	50.0	33.3	66.7	37.1
Claude-3.5	2	32.0	25.0	50.0	33.3	66.7	30.6
	3	30.0	75.0	0.0	66.7	0.0	32.3
	1	38.6	0.0	100.0	11.1	100.0	39.1
OpenAI-o1	2	21.1	0.0	0.0	0.0	0.0	18.8
	3	40.4	100.0	0.0	0.0	0.0	42.2
	1	41.7	33.3	100.0	25.0	40.0	41.9
GPT-4o	2	31.7	0.0	0.0	0.0	0.0	25.7
	3	26.7	66.7	0.0	75.0	60.0	32.4

Table 4: **Correction Rate Results across various LMMs on InterFeedback-Human.** Math<sup>Text</sup> and Coding<sup>Text</sup> represent two text-only task categories. # Round denotes the number of interaction rounds. The correction rate is the percentage of corrected samples among all erroneous samples.

The best scores for each subcategory in our InterFeedback-Human are 37.5% (Claude-3.5-Sonnet), 70.0% (GPT-40), 90% (OpenAI-01), and 90% (OpenAI-01), respectively. (2) Overall, Claude-3.5 achieves the highest average accuracy at 48.3%.

Correction rate results analysis. Comparing the correction rates across rounds in Table 4 reveals that GPT-40 benefits the most from human feedback in the first round, correcting 41.9% of erroneous samples, while Claude-3.5 exhibits its strongest correction performance in the second round, with 30.6% of erroneous samples corrected. Given that the ground truth answer is provided in the third round, all LMMs are able to supply their reasoning steps for selecting the correct answer.

**Distribution of Tasks Corrected Across Rounds.** Figure 3 illustrates the distribution of tasks solved by each LMM across the interaction rounds. Round 0 represents the initial accuracy before beginning human-AI interactions. For example, GPT-40 solved 38.3% of instances in Round 0, 25.8% in Round 1, and 20% in Round 2. Additionally, during the first two rounds, both OpenAI-o1 and Claude-3.5-Sonnet solved the same number of

samples, achieving a performance of 67.5%.

**Distribution of corrected samples across various task categories.** As shown in Figure 4, Visual logic tasks are mostly resolved within the first two rounds, whereas Math (Text-only) and MMMU-Pro tasks show little corrections in rounds 1 and 2. In contrast, Coding (Text-only) and MathVerse tasks exhibit corrections during rounds 1 and 2.

#### **Conclusion**

In this work, we introduced InterFeedback-Bench, the first solution to concern the critical importance of evaluating the interactive intelligence of current LMMs. We build an interactive framework InterFeedback which can be applied to any LMM and dataset to bootstrap the testing in an interactive way. We conduct the comprehensive evaluations on 10 open-source LMMs by demonstrating with two representative datasets MathVerse and MMMU-Pro. Additionally, we present InterFeedback-Human, a new benchmark for manually testing the proprietary models such as OpenAI-o1 and Claude-3.5 with 120 curated samples. Our evaluation results show that even the SOTA LMM (like OpenAI-o1) can only correct their results through human feed-

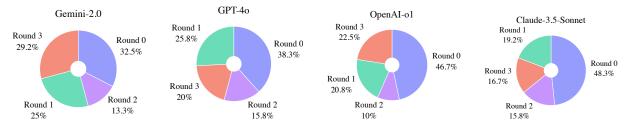


Figure 3: Distribution of samples being corrected in each round. We can observe that Claude-3.5-Sonnet archives the best performance in round 0.

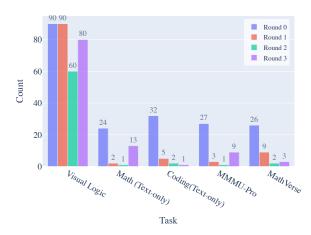


Figure 4: Distribution of corrected samples across various task categories. Visual logic tasks are mostly resolved within the first two rounds, whereas Math (Textonly) and MMMU-Pro tasks show few corrections in rounds 1 and 2. In contrast, Coding (Text-only) and MathVerse tasks exhibit corrections during rounds 1 and 2.

back with less than 50%. Several findings point to the essential need for methods that improve the LMM's ability to receive feedback to improve itself.

#### Limitations

Our method is not without limitations. First, as an initial attempt in the field, this work proposes a straightforward method to bootstrap the LMMs in an interactive way. We use the proprietary LMM to stimulate the humans mimicking the human-AI interaction process. Due to the difficulty of existing benchmarks, the proprietary LMMs may not fully provide all pertinent feedback, though we propose two strategies: 1) select the intersection set for testing and 2) record the valid output only. Due to the limitation of GPU memory, we have to select the tested LMMs within 7B parameters.

# Acknowledgement

This project is supported by Mike Zheng Shou's Start-Up Grant (A-0009453-03-00).

### References

Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi, Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. Phi-3 technical report: A highly capable language model locally on your phone. Preprint, arXiv:2404.14219.

Anthropic. 2024. Introducing computer use, a new claude 3.5 sonnet, and claude 3.5 haiku.

- Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşırlar. 2023. Introducing our multimodal models.
- Sanxing Chen, Sam Wiseman, and Bhuwan Dhingra. 2024a. Chatshop: Interactive information seeking with language agents. *Preprint*, arXiv:2404.09911.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. 2015. Microsoft coco captions: Data collection and evaluation server. *arXiv:1504.00325*.
- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 2024b. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24185–24198.
- Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert, Yvonne Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat, Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Borchardt, Dirk Groeneveld, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick, Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick, Ali Farhadi, and Aniruddha Kembhavi. Molmo and pixmo: Open weights and open data for state-of-the-art vision-language models. Preprint, arXiv:2409.17146.
- Lakshita Dodeja, Pradyumna Tambwekar, Erin Hedlund-Botti, and Matthew Gombolay. 2024. Towards the design of user-centric strategy recommendation systems for collaborative human—ai tasks. *International Journal of Human-Computer Studies*, 184:103216.
- Difei Gao, Lei Ji, Zechen Bai, Mingyu Ouyang, Peiran Li, Dongxing Mao, Qinchen Wu, Weichen Zhang, Peiyi Wang, Xiangwu Guo, Hengxu Wang, Luowei Zhou, and Mike Zheng Shou. 2024. Assistgui: Taskoriented pc graphical user interface automation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13289–13298.
- Gemini. 2024. Our next-generation model: Gemini 1.5.

Gemini. 2025. Gemini 2.0.

Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*.

- Imbue. 2024. Imbue code comprehension.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*.
- Akbir Khan, John Hughes, Dan Valentine, Laura Ruis, Kshitij Sachan, Ansh Radhakrishnan, Edward Grefenstette, Samuel R. Bowman, Tim Rocktäschel, and Ethan Perez. 2024. Debating with more persuasive llms leads to more truthful answers. *Preprint*, arXiv:2402.06782.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024a. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*.
- Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. 2023. Seed-bench: Benchmarking multimodal llms with generative comprehension. *Preprint*, arXiv:2307.16125.
- Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. 2024b. Numinamath. [https://huggingface.co/AI-MO/NuminaMath-CoT](https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina\_dataset.pdf).
- Lei Li, Yuancheng Wei, Zhihui Xie, Xuqing Yang, Yifan Song, Peiyi Wang, Chenxin An, Tianyu Liu, Sujian Li, Bill Yuchen Lin, Lingpeng Kong, and Qi Liu. 2024c. Vlrewardbench: A challenging benchmark for vision-language generative reward models. *Preprint*, arXiv:2411.17451.
- Pengxiang Li, Zhi Gao, Bofei Zhang, Tao Yuan, Yuwei Wu, Mehrtash Harandi, Yunde Jia, Song-Chun Zhu, and Qing Li. 2024d. FIRE: A dataset for feedback integration and refinement evaluation of multimodal models. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*.
- Shuyue Stella Li, Vidhisha Balachandran, Shangbin Feng, Jonathan S. Ilgen, Emma Pierson, Pang Wei Koh, and Yulia Tsvetkov. 2024e. Mediq: Questionasking LLMs and a benchmark for reliable interactive clinical reasoning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Kevin Qinghong Lin, Linjie Li, Difei Gao, Qinchen Wu, Mingyi Yan, Zhengyuan Yang, Lijuan Wang, and Mike Zheng Shou. 2024. Videogui: A benchmark for gui automation from instructional videos. *arXiv* preprint arXiv:2406.10227.

- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2024a. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024b. Llavanext: Improved reasoning, ocr, and world knowledge.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023a. Visual instruction tuning. In *NeurIPS*.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. 2023b. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2024. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *International Conference on Learning Representations* (*ICLR*).
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. *NeurIPS*.
- Ziyang Luo, Haoning Wu, Dongxu Li, Jing Ma, Mohan Kankanhalli, and Junnan Li. 2024. Videoautoarena: An automated arena for evaluating large multimodal models in video analysis through user simulation. *Preprint*, arXiv:2411.13281.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *CVPR*.
- Scott Mayer McKinney, Marcin Sieniek, Varun Godbole, Jonathan Godwin, Natasha Antropova, Hutan Ashrafian, Trevor Back, Mary Chesus, Greg S Corrado, Ara Darzi, et al. 2020. International evaluation of an ai system for breast cancer screening. *Nature*, 577(7788):89–94.
- OpenAI. 2023. Gpt-4o.
- OpenAI. 2024. Openai o1 system card.
- OpenGVLab. 2024. Internvl2: Better than the best—expanding performance boundaries of open-source multimodal models with the progressive scaling strategy.
- Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22.

- Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022. A-okvqa: A benchmark for visual question answering using world knowledge. *arXiv*.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards vqa models that can read. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- ByteDance Seed Team. 2025a. Seed1.5-vl technical report. *arXiv preprint arXiv*:2505.07062.
- Qwen Team. 2025b. Qwen2.5-vl.
- Maria Virvou. 2022. The emerging era of human-ai interaction: Keynote address. In 2022 13th International Conference on Information, Intelligence, Systems & Applications (IISA), pages 1–10. IEEE.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *Preprint*, arXiv:2409.12191.
- Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. 2023. Cogvlm: Visual expert for pretrained language models. *Preprint*, arXiv:2311.03079.
- Haoran Wei, Youyang Yin, Yumeng Li, Jia Wang, Liang Zhao, Jianjian Sun, Zheng Ge, Xiangyu Zhang, and Daxin Jiang. 2025. Slow perception: Let's perceive geometric figures step-by-step. *Preprint*, arXiv:2412.20631.
- Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhou-jun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. 2024. Osworld: Benchmarking multimodal agents for openended tasks in real computer environments. *Preprint*, arXiv:2404.07972.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing

- Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. 2024. Qwen2 technical report. *Preprint*, arXiv:2407.10671.
- John Yang, Carlos E Jimenez, Alex L Zhang, Kilian Lieret, Joyce Yang, Xindi Wu, Ori Press, Niklas Muennighoff, Gabriel Synnaeve, Karthik R Narasimhan, Diyi Yang, Sida Wang, and Ofir Press. 2025. SWE-bench multimodal: Do autonomous programming systems generalize to new software domains? In *The Thirteenth International Conference on Learning Representations*.
- John Yang, Akshara Prabhakar, Karthik R Narasimhan, and Shunyu Yao. 2023. Intercode: Standardizing and benchmarking interactive coding with execution feedback. In Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track.
- Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik R Narasimhan. 2025. {\$\tau\$}-bench: A benchmark for \underline{T}ool-\underline{A}gent-\underline{U}ser interaction in real-world domains. In *The Thirteenth International Conference on Learning Representations*.
- Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, Qianyu Chen, Huarong Zhou, Zhensheng Zou, Haoye Zhang, Shengding Hu, Zhi Zheng, Jie Zhou, Jie Cai, Xu Han, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. Minicpmv: A gpt-4v level mllm on your phone. *Preprint*, arXiv:2408.01800.
- Se-eun Yoon, Zhankui He, Jessica Echterhoff, and Julian McAuley. 2024. Evaluating large language models as generative user simulators for conversational recommendation. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 1490–1504, Mexico City, Mexico. Association for Computational Linguistics.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. 2024a. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567.
- Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, Yu Su, Wenhu Chen, and Graham Neubig. 2024b. Mmmu-pro: A more robust

- multi-discipline multimodal understanding benchmark. *Preprint*, arXiv:2409.02813.
- Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6720–6731.
- Jiehuang Zhang, Ying Shu, and Han Yu. 2021. Human-machine interaction for autonomous vehicles: A review. In *International Conference on Human-Computer Interaction*, pages 190–201.
- Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Peng Gao, and Hongsheng Li. 2024. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? *Preprint*, arXiv:2403.14624.
- Tianyi Zhang, Isaac Tham, Zhaoyi Hou, Jiaxuan Ren, Liyang Zhou, Hainiu Xu, Li Zhang, Lara J Martin, Rotem Dror, Sha Li, et al. 2023. Human-in-the-loop schema induction. *arXiv:2302.13048*.
- Hengyuan Zhao, Pan Zhou, Difei Gao, Zechen Bai, and Mike Zheng Shou. 2024a. LOVA3: Learning to visual question answering, asking and assessment. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Henry Hengyuan Zhao, Difei Gao, and Mike Zheng Shou. 2025. Worldgui: Dynamic testing for comprehensive desktop gui automation. *Preprint*, arXiv:2502.08047.
- Henry Hengyuan Zhao, Pan Zhou, and Mike Zheng Shou. 2024b. Genixer: Empowering multimodal large language model as a powerful data generator. In *European Conference on Computer Vision*, pages 129–147. Springer.

#### A Model Sources.

For different LMMs, we select their latest models with sizes around 7B for evaluation. Table 5 presents the release time and model sources of LMMs used in InterFeedback-Bench.

## **B** Error Analysis.

The Table 6 shows the answers generated after receiving feedback from two feedback providers, along with the corrected flag compared to the GT answer. Across these 10 questions, 11 out of 20 samples could not be corrected. Half of the questions remained uncorrected regardless of the feedback provider. Except for the first question in the table, the LMM generates identical answers after receiving feedback from different providers. This pattern implies that the model's capacity to improve based on feedback resembles an inherent ability. Prompting LMM with feedback can be considered another prompting strategy to invoke the inherent ability rather than a robust ability to reason and incorporate new information to address challenging questions.

# C Qualitative Analysis.

Interactive process could improve the performance of proprietary LMMs. In Figure 5, we provide the qualitative results of different models. For the same question, Claude-3.5-Sonnet gives the correct answer C without human feedback, Gemini-2.0-Flash uses two rounds while OpenAI-o1 uses three rounds. It indicates that 1) even the SOTA models like OpenAI-o1 can not fully address the visual logic problem which is worse than Claude-3.5-Sonnet, 2) the responses can be corrected by human feedback which shows that the models have the capability of interpreting and incorporating the feedback into their reasoning, 3) Different models shows a different level of this capability. Additionally, we provide another example in Figure 6.

LMMs may not truly reasoning-They guess answers by elimination. In Figure 7, we find that the model will guess the answer when we only have four options, the model tends to guess answers. For the same question, we conduct twice runs and find that OpenAI-o1 could not solve this problem at the beginning, but two different answers were given in these two runs. In the first run, the model outputs D at the beginning while in the second run, the model outputs the A at the beginning. In the following rounds, we provide the same prompts to ensure the

fairness comparison, one can see that based on the same prompt, it outputs the same answer C in the second round. The left run in the figure shows the correct answer in the third round while the right run in the figure shows the incorrect answer D. We continue to give the third feedback for round 4, and the right run finally gives answer B. It is obvious that when a problem cannot solved by a model, it will 1) outcome answer randomly, and 2) outcome the answer through an elimination approach. These results may indicate that LMMs may not always truly reason they may give the answer by guessing. Additionally, we provide another example in Figure 8 to illustrate that LMMs may guess answers when they can not solve the challenging problems.

LMMs still fail when the GT answer is not provided in the level 3 feedback. As discussed in Section 3.2, we include the GT answer in the level 3 feedback prompt to examine whether the model can generate the correct reasoning procedure that leads to the correct answer. When we remove the GT answer as in Figure 9, the model still fails to produce the correct answer, indicating its limited capability in solving challenging problems even when detailed feedback is provided as guidance.

# D Examples of Feedback.

We provide the examples of feedback provided by Claude-3.5-Sonnet on MathVerse and MMMU-Pro, respectively. As these examples show that after providing the feedback, the questions are solved correctly, and the provided feedback is concise and pertinent without leaking the GT. The feedback are mainly focused on analyzing the question and providing the reasoning thoughts for the tested LMMs to use the additional information for solving questions.

prompt\_system= """You are a reasoning assistant tasked with solving questions based on visual analysis. Your goal is to analyze the provided question and image carefully and provide the correct answer. Instructions: 1. Focus on observing all relevant details in the image. 2. Ensure thorough and logical reasoning before arriving at a conclusion. 3. Respond concisely by selecting only a single letter: A, B, C, or D. Remember, no explanations or additional text are needed—only the chosen letter as the answer."""

Model	odel Release Time Source					
		Proprietary Models				
GPT-40 (OpenAI, 2023)	2024-08-26	https://openai.com/index/hello-gpt-4o/				
OpenAI-o1 (OpenAI, 2024)	2024-12-17	https://openai.com/o1/				
Gemini-1.5-Flash (Gemini, 2024)	2024-09-24	https://deepmind.google/technologies/gemini/				
Gemini-2.0-Flash	2025-01-21	https://deepmind.google/technologies/gemini/				
Claude-3.5-Sonnet	2024-10-22	https://www.anthropic.com/claude/sonnet				
Claude-Sonnet-4	2025-05-23	https://www.anthropic.com/news/claude-4				
Open-source Models						
LLaVA-One-Vision	2024-08-05	https://llava-vl.github.io/blog/2024-08-05-llava-onevision/				
InterVL2-8B	2024-07-04	https://internvl.github.io/blog/2024-07-02-InternVL-2.0				
InterVL3-8B	2025-04-11	https://huggingface.co/OpenGVLab/InternVL3-8B				
Molmo-7B	2024-09-24	https://huggingface.co/allenai/Molmo-7B-D-0924				
MiniCPM-V	2024-08-03	https://huggingface.co/openbmb/MiniCPM-V				
GLM-4V-9B	2024-11-01	https://huggingface.co/THUDM/glm-4v-9b				
Pih3.5-Vision-4.2B	2024-08-20	https://huggingface.co/microsoft/Phi-3.5-vision-instruct				
LLaVA-1.5-7B	2023-10-05	https://huggingface.co/liuhaotian/llava-v1.5-7b				
LLaVA-1.6-Mistral-7B	2024-01-30	https://huggingface.co/llava-hf/llava-v1.6-mistral-7b-hf				
Fuyu-8B	2023-10-27	https://huggingface.co/adept/fuyu-8b				
Qwen2-VL-7B	2024-08-30	https://huggingface.co/Qwen/Qwen2.5-VL-7B-Instruct				
Qwen2.5-VL-7B	2025-02-20	https://huggingface.co/Qwen/Qwen2.5-VL-7B-Instruct				
Seed-1.5-VL-Thinking	2025-05-13	https://github.com/ByteDance-Seed/Seed1.5-VL				

Table 5: The release time and model source of LMMs used in our InterFeedback-Bench.

Task ID	Enhanced by Gemini-1.5-Flash	Has it been corrected?	Enhanced by Claude-3.5-Sonnet	GT Answer	Has it been corrected?
History_134	В	Yes	A	В	No
Pharmacy_287	A	No	A	В	No
validation_Sociology_29	C	No	C	В	No
Literature_19	D	No	D	В	No
Agriculture_169	C	Yes	C	C	Yes
validation_Economics_12	В	No	В	A	No
Electronics_68	A	Yes	A	A	Yes
Design_126	D	Yes	D	D	Yes
Medical_Science_28	A	Yes	A	A	Yes
test Pharmacy 140	C	No	C	В	No

Table 6: Comparison of Enhancements and Corrections on cases from MMMU-Pro(Yue et al., 2024b).

prompt\_template = """ Please carefully analyze the image and choose the correct answer to the following question.

{question}

Take your time to observe all relevant details in the image and ensure your reasoning is thorough before selecting your answer. Answer with only a single letter (A, B, C, or D): """

prompt\_template\_local = """ You are a reasoning assistant tasked with solving questions based on visual analysis. Your goal is to analyze the question and image carefully and provide the correct answer. Instructions:

- 1. Focus on observing all relevant details in the image.
- 2. Ensure thorough and logical reasoning before arriving at a conclusion. {question}

Answer with only a single letter (A, B, C, or D): """

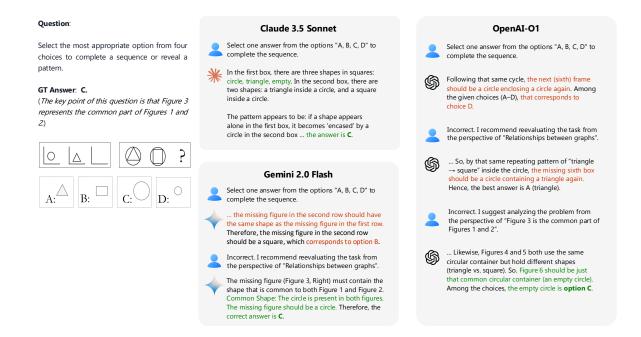


Figure 5: Qualitative results on different LMMs.

prompt\_feedback\_system = You are tasked with generating constructive feedback when solving questions.

**Instructions:** 

- 1. Analyze the provided question and the incorrect answer.
- 2. Identify general reasoning flaws or oversights that might have led to the incorrect response.
- 3. Suggest how the reasoning can be improved, such as paying more attention to specific observations, question details, or logical deductions.
- 4. DO NOT mention the correct answer or any specific details of the question.

Your feedback should be constructive, and helpful.

prompt with feedback = """ You are an assistant tasked with re-evaluating a question after receiving feedback on a previous incorrect answer.

Instructions: 1. Carefully review the question and the provided feedback.

- 2. Use the feedback to refine your reasoning and identify details or logical steps you may have overlooked.
- 3. Re-examine all aspects of the question thoroughly before selecting your answer. question

Your previously selected answer: {predicted answer} was incorrect.

Feedback: {feedback}

Answer with only a single letter (A, B, C, or D): """

prompt\_feedback = Question: {question} Previous answer: '{predicted\_answer}' was

Provide a feedback that helps correctly answer the question.

incorrect.

#### Ouestion: Claude 3.5 Sonnet OpenAI-01 As shown in the figure, it is known that angle Select one answer from the options "A, B, C, D" to Select one answer from the options "A, B, C, D" to A = 80°, angle B = 60, DE parallel BC, then complete the sequence. the size of angle CED is \_\_\_? ... Since we know $\angle$ A = 80° and $\angle$ B = 60 ° , We can find $\angle$ C in triangle ABC: $\angle$ C = 180° - $\angle$ A - $\angle$ B, Since DE || BC, triangle ADE is similar to triangle Choices: A: 40° B: 60° C: 120° D: 140° $\angle C = 180^{\circ} - 80^{\circ} - 60^{\circ} = 40^{\circ}.$ Line AC is a transversal of the parallel lines DE Now, because DE is parallel to BC, $\angle$ CED = 180° - $\angle$ C, $\angle$ CED = 180° - 40° = 140°. GT Answer: D: 140°. and BC, so the angle formed at E by CE and DB (i.e. $\angle$ CED) corresponds to $\angle$ C in triangle ABC. Therefore, angle CED is 140°. The answer is **D: 140°.** In triangle ABC, the angles are 80° (at A), 60° (at B), and thus 40° (at C). Therefore, ∠CED = 40°. Gemini 2.0 Flash D Incorrect, I recommend reevaluating the task from Select one answer from the options "A, B, C, D" to the perspective of thinking about the triangle as complete the sequence. 180° and the parallel lines' properties В Since DE is parallel to BC, the corresponding ... The angles ∠BCE (at C) and ∠CED (at E) lie on the angles are equal. Angle CED is a corresponding angle to angle C. ..., Therefore, ... The correct answer is A: 40°. same side of this transversal, so they are co-interior angles and must sum to 180°. Incorrect. I recommend reevaluating the task from the perspective of thinking about the triangle as $\angle$ BCE is the same as $\angle$ BCA = 40°. Hence, $\angle$ CED = 180° - 40° = **140°**.

Figure 6: Qualitative results on different LMMs.

180° and the parallel lines' properties.

Since DE is parallel to BC, the consecutive interior angles on the same side of the transversal are supplementary. Therefore, ... The correct answer is D: 140°.

prompt\_with\_simple\_feedback = """ You
are an assistant tasked with re-evaluating a
question on a previous incorrect answer.
Instructions:

- 1. Carefully review the question.
- 2. Refine your reasoning and identify details or logical steps you may have overlooked.
- 3. Re-examine all aspects of the question thoroughly before selecting your answer. {question}

Your previously selected answer: {predicted\_answer} was incorrect.

Answer with only a single letter (A, B, C, or D): """

ID	Initial Answer	New Answer(simple)	New Answer(Detail)	GT	Feedback
73	A	В	D	D	"Let me provide some constructive feedback to help with this geometry problem:  1. When analyzing angles formed by intersecting lines, remember these key principles: - Adjacent angles sum to 180° - When one line rotates, it affects multiple angles simultaneously - The rotation angle is related to the difference between initial and final positions  2. Common mistakes to avoid: - Don't assume the rotation angle is the same as the given angle - Consider what the final position needs to be (parallel to line a) - Think about the smallest rotation needed to achieve the desired result  3. Suggested approach: - First identify the target position (where line b needs to end up) - Look at the relationship between angle 1 and the amount of rotation needed - Consider that the rotation should make the lines parallel - Calculate how much movement is needed from the current position to the final position
102	A	В	D	D	"Here's some constructive feedback to help solve this problem:  1. When analyzing angles formed by intersecting lines, remember to look for:  - Corresponding angles - Alternate interior angles - Alternate exterior angles - Supplementary angles (angles that add up to 180°)  2. In this case, pay special attention to:  - The relationship between angle 4 and the other angles - The fact that angle 1 and angle 2 are on the same side of a line - How these angles relate to form a complete picture  3. Key points to improve your approach:  - Draw out all the angle relationships clearly - Remember that when two parallel lines are cut by a transversal, certain angle relationships are formed - Use the given information that angles 1, 2, and 3 are all equal to 55°"

Table 7: Feedback examples provided by Claude-3.5-Sonnet on MathVerse Dataset.

ID	Initial Answer	New Answer(simple)	New Answer(Detail)	GT	Feedback
test_Pharmacy_287	A	B	B	В	"Here's some constructive feedback to help with answering this type of question:  1. When analyzing protein structure questions, carefully examine the image for key structural characteristics: Look for distinct individual polypeptide chains/subunits - Notice how the helices are arranged and connected - Consider whether the structure appears to be a single unit or multiple units working together  2. Remember the definition of quaternary structure: - Review what specifically defines quaternary structure in proteins - Don't confuse complex folding patterns within a single chain with quaternary structure  3. Methodology for improvement: First identify all visible protein chains in the diagram - Look for points of interaction between separate chains - Don't assume that just because a protein structure looks complex, it automatically has
test_Math_270	C	A	A	A	quaternary structure "Here's some constructive feedback to help solve this problem:  1. Approach Organization: - Before diving into calculations, take a moment to visualize what's happening: a lad der sliding down a wall with its bottom moving outward at a constant rate - Consider what geometric principles are in volved (right triangle, Pythagorean the orem) - Think about what information is specifically asked for (rate of change of height)  2. Key Considerations: - Pay attention to the specific moment you're asked about (when the top reaches the wall height) - Remember that this is a related rates problem where one rate is given (bottom moving out) and another is asked for (top moving down) - Consider drawing the triangle at the specific instant mentioned in the question  3. Problem-Solving Strategy: - Set up your equation carefully before plugging in numbers - When dealing with rates of change, differentiation will be needed

Table 8: Feedback examples provided by Claude-3.5-Sonnet on MMMU-Pro Dataset.

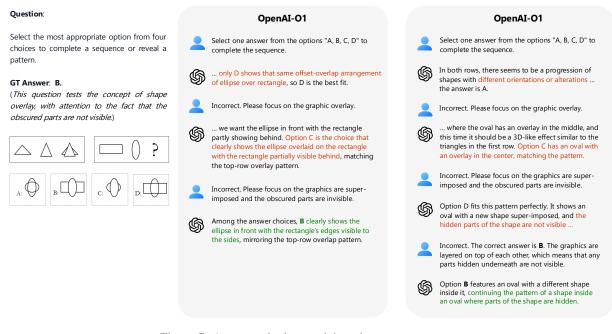


Figure 7: An example that model tends to guess answers.

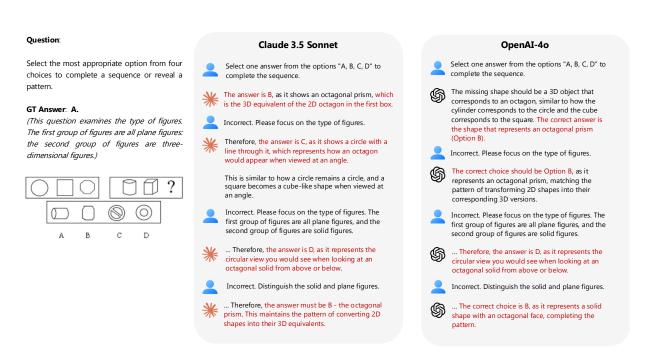


Figure 8: An example that model tends to guess answers.

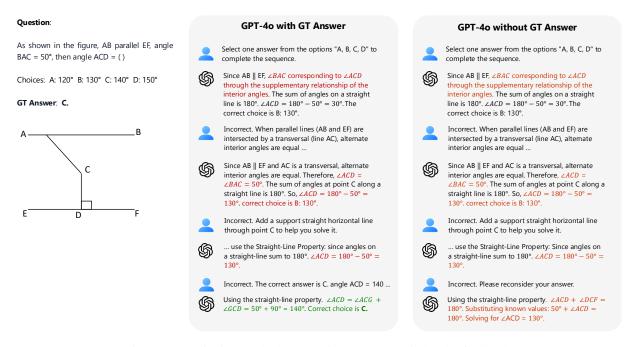


Figure 9: Qualitative results by removing GT answer in level 3 feedback.