## MMEVALPRO: Calibrating Multimodal Benchmarks Towards Trustworthy and Efficient Evaluation

Jinsheng Huang<sup>1,2,3\*</sup>, Liang Chen<sup>1, 2\*</sup>, Taian Guo<sup>1,2,3</sup>, Fu Zeng<sup>4</sup>, Yusheng Zhao<sup>1,2,3</sup> Bohan Wu<sup>1,2,3</sup>, Ye Yuan<sup>1,2,3</sup>, Haozhe Zhao<sup>1</sup>, Zhihui Guo<sup>5</sup>, Yichi Zhang<sup>1</sup>, Jingyang Yuan<sup>1,2,3</sup> Wei Ju<sup>1,2,3</sup>, Luchen Liu<sup>1,2,3</sup>, Tianyu Liu<sup>6</sup>, Baobao Chang<sup>1, 2†</sup>, Ming Zhang<sup>1, 2, 3†</sup> <sup>1</sup>National Key Laboratory for Multimedia Information Processing, Peking University

<sup>2</sup>School of Computer Science, Peking University, <sup>3</sup>PKU-Anker LLM Lab

<sup>4</sup>Chinese Academy of Medical Sciences, <sup>5</sup>CUHK, <sup>6</sup>Alibaba Group

## Abstract

Large Multimodal Models (LMMs) exhibit impressive cross-modal understanding and reasoning abilities, often assessed through multiplechoice questions (MCQs) that include an image, a question, and several options. However, many benchmarks used for such evaluations suffer from systematic biases. Remarkably, Large Language Models (LLMs) without any visual perception capabilities achieve nontrivial performance, undermining the credibility of these evaluations. To address this issue while maintaining the efficiency of MCQ evaluations, we propose MMEVALPRO, a benchmark designed to avoid Type-I errors through a trilogy evaluation pipeline and more rigorous metrics. For each original question from existing benchmarks, human annotators augment it by creating one perception question and one knowledge anchor question through a meticulous annotation process. MMEVALPRO comprises 2,138 question triplets, totaling 6,414 distinct questions. Two-thirds of these questions are manually labeled by human experts, while the rest are sourced from existing benchmarks (MMMU, ScienceQA, and MathVista). Compared with the existing benchmarks, our experiments with the latest LLMs and LMMs demonstrate that MMEVALPRO is more challenging (the best LMM lags behind human performance by 31.73%, compared to an average gap of 8.03% in previous benchmarks) and more trustworthy (the best LLM trails the best LMM by 23.09%, whereas the gap for previous benchmarks is just 14.64%). Our indepth analysis explains the reason for the large performance gap and justifies the trustworthiness of evaluation, underscoring its significant potential for advancing future research.

## 1 Introduction

Ever since the birth of standardized testing, the credibility of its conclusions has been a signif-



Figure 1: Examples of the probing experiments.



Figure 2: Topic distribution of MMEVALPRO's data.

icant concern. The same problem goes for the evaluation of recently popular Large Multimodal Models (LMMs) such as GPT4-0 (OpenAI, 2024), Gemini-1.5 (Team et al., 2024), Qwen-VL (Bai et al., 2023b) and LLaVA (Liu et al., 2023b). One classic composition of such an evaluation is the multiple-choice question (MCQ), which includes an image, a question, possible choices, and an answer. This form of evaluation has higher usability

<sup>&</sup>lt;sup>\*</sup>Equal contribution. <sup>†</sup>Corresponding authors.