

Large Multimodal Model Robustness to Distributional Shifts in Chart-Based Tasks

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How robust are LMMs trained on MMC-Instruction to distributional shifts in chart types or domains, as quantified by cross-domain accuracy when tested on unseen chart datasets. 12 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.9/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: MMC: Advancing Multimodal Chart Understanding with Large-scale Instruction Tuning. Research question: How robust are LMMs trained on MMC-Instruction to distributional shifts in chart types or domains, as quantified by cross-domain accuracy when tested on unseen chart datasets?.

2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.9/10.

3 Results

13 papers retrieved. 12 claims extracted; 2 independently verified. Quality review score: 5.9/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
MMC-Benchmark is a human-annotated benchmark containing nine distinct tasks for evaluating reasoning capabilities over c	✓	0.25
Existing Large Multimodal Models (LMMs), including GPT-4V, show limitations in correctly interpreting charts when evalua	✓	0.22
GPT-4V faces significant challenges on MMC-Benchmark, specifically in 'Chart to Datatable' and 'Chart to Json' tasks.	×	0.10
The MMC-Instruction dataset contains 600k samples, making it larger than FigureQA (180k), DVQA (300k), PlotQA (224k), Ch	×	0.06
MMC-Instruction supports free-form answers and open-ended/MQA formats, whereas FigureQA, DVQA, and PlotQA use fixed voca	×	0.04
MMCA achieves state-of-the-art performance on current chart question-answer benchmarks compared with existing open-sourc	×	0.12
On the ChartQA benchmark, MMCA achieves a score of 57.4, outperforming the variant without fine-tuned vision encoder whi	×	0.03
On the DocVQA benchmark, MMCA achieves a score of 72.5, outperforming the variant without fine-tuned vision encoder whic	×	0.03
On the TextVQA benchmark, MMCA achieves a score of 59.6, outperforming the variant without fine-tuned vision encoder whi	×	0.03
Fine-tuning the vision encoder part of the MMCA model is necessary for optimal performance, as the model under-performs	×	0.06
In a specific test case regarding land area, GPT-4V and LLaVA-v1.5 incorrectly identified China as the third largest cou	×	0.03
MMC-Benchmark includes tasks such as chart information extraction, chart reasoning, contextual chart understanding, char	×	0.13

References

- <http://arxiv.org/abs/2312.05435v1>

- <http://arxiv.org/abs/2311.10774v2>
- <http://arxiv.org/abs/2407.14506v3>