

Self-Repair Transferability in Llama-2 Across Diverse Instruction-Tuning Domains

Assignee Research

June 2, 2026

Abstract

This report synthesises findings from 7 peer-reviewed papers addressing the following research question: How does the cross-domain transferability of self-repair mechanisms in Llama-2 models scale with instruction-tuning data diversity, as measured by pass@k accuracy across different programming. 10 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Vision-Flan: Scaling Human-Labeled Tasks in Visual Instruction Tuning. Research question: How does the cross-domain transferability of self-repair mechanisms in Llama-2 models scale with instruction-tuning data diversity, as measured by pass@k accuracy across different programming language domains?.

2 Methodology

Systematic literature search across multiple databases yielded 7 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.2/10.

3 Results

7 papers retrieved. 10 claims extracted; 1 independently verified. Quality review score: 4.2/10.

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
VISION-FLAN BASE achieves state-of-the-art performance on comprehensive evaluation benchmarks including MME, MM-Bench and	×	0.11
VISION-FLAN BASE reduces hallucination and catastrophic forgetting.	×	0.11
VISION-FLAN BASE scores significantly lower on the LLaVA-Bench dataset compared to VLMs trained using GPT-4 synthesized	×	0.14
VISION-FLAN CHAT achieves significant performance improvement on LLaVA-Bench through the second-stage tuning on a mere 1	✓	0.18
Training on academic datasets leads VISION-FLAN BASE to generate brief responses, which are not aligned with human preference	×	0.09
Visual instruction tuning mainly enables LLMs to better understand visual features while MLPs have been sufficiently learned	×	0.14
Replacing the instruction-tuned MLPs in VISION-FLAN BASE and VISION-FLAN CHAT with the pre-trained MLPs from the pre-training	×	0.06
The Pearson Correlation Coefficient between the parameters of pretrained MLPs and instruction-tuned MLPs is computed.	×	0.03
VISION-FLAN dataset contains 1.6M instances and 196 tasks.	×	0.06
VISION-FLAN dataset is based on publicly available datasets.	×	0.10

References

- <http://arxiv.org/abs/2402.11690v1>

- <http://arxiv.org/abs/2306.09896v5>
- <http://arxiv.org/abs/2312.10793v3>