

Comparative Analysis of Entity-Centric Multimodal Preference Optimization and Pareto-Based Alignment for Hallucination Reduction

Assignee Research

June 13, 2026

Abstract

Large Visual Language Models (LVLMs) have demonstrated impressive capabilities across multiple tasks. However, their trustworthiness is often challenged by hallucinations, which can be attributed to the modality misalignment and the inherent hallucinations of their underlying Large Language Models (LLMs) backbone. Existing preference alignment methods focus on aligning model responses with human preferences while neglecting image-text modality alignment, resulting in over-reliance on LLMs and hallucinations. In this paper, we propose Entity-centric Multimodal Preference Optimization (EMPO), wh

1 Introduction

This paper examines: Mitigating Hallucinations in Large Vision-Language Models via Entity-Centric Multimodal Preference Optimization. Research question: How does entity-centric multimodal preference optimization compare to traditional Pareto-based alignment in reducing hallucination rates on the MMHal-Bench dataset when evaluated across different visual complexity levels?.

2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.2/10.

3 Results

15 papers retrieved. 16 claims extracted; 10 independently verified. Quality review score: 7.2/10.

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
EMPO enhances multimodal semantic alignment across five widely-used benchmarks.	✓	0.16
EMPO effectively reduces hallucinations across five widely-used benchmarks.	✓	0.17
Recent research constructs Large Vision-Language Models (LVLMs) by aligning Large Language Models (LLMs) with visual mod	×	0.13
Recent LVLMs demonstrate superior performance across various visual-language tasks compared to earlier studies.	×	0.15
LVLMs typically adopt a two-stage training strategy consisting of pretraining on large-scale image-text pairs and instru	✓	0.22
LLaVA introduces synthetic instructions to fine-tune an instruction-following LVLM.	✓	0.18
MiniGPT-v2 employs unique task identifiers during fine-tuning to reduce instruction ambiguity.	✓	0.19
Hallucinations in LVLMs occur when model responses conflict with images, instructions, or context.	×	0.10
Some methods mitigate hallucinations by filtering out long-tail or entity co-occurrence data.	×	0.11
Filtering out long-tail or entity co-occurrence data involves high annotation costs.	✓	0.16
Post-processing techniques such as optimizing decoding strategies or applying post-hoc corrections reduce hallucinations	✓	0.18
LLaVA-RLHF pioneered the exploration of human preference alignment to mitigate hallucinations in LVLMs.	×	0.13
RLHF-V, RLAIIF-V, and POVID refined human preference alignment approaches with improved visual localization, text segment	✓	0.16
Existing human preference alignment methods primarily focus on response-level preferences while neglecting the requireme	✓	0.21
MDPO proposed image-conditional preference alignment but overlooked aligning instructions with human preferences.	✓	0.18
EMPO incorporates preferences across comprehensive aspects.	×	0.11

References

- <http://arxiv.org/abs/2503.14504v2>
- <http://arxiv.org/abs/2506.04039v2>
- <http://arxiv.org/abs/2506.11712v3>