

Visual Reasoning Transferability in Large Multimodal Models Across Domains

Assignee Research

June 6, 2026

Abstract

This report synthesises findings from 9 peer-reviewed papers addressing the following research question: What is the cross-domain transferability of visual reasoning capabilities in LMMs when trained on HumanEval-V versus traditional multimodal benchmarks like VQA or COCO. 15 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Advances in Multimodal Adaptation and Generalization: From Traditional Approaches to Foundation Models. Research question: What is the cross-domain transferability of visual reasoning capabilities in LMMs when trained on HumanEval-V versus traditional multimodal benchmarks like VQA or COCO?.

2 Methodology

Systematic literature search across multiple databases yielded 9 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.0/10.

3 Results

9 papers retrieved. 15 claims extracted; 0 independently verified. Quality review score: 3.0/10.

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
Zanella et al. [122] improve zero-shot and few-shot VLM performance by jointly optimizing view quality assessment and de	×	0.01
Zhang et al. [123] introduce dual memory networks, where static memory caches training data knowledge and dynamic memory	×	0.02
Wang et al. [124] applied Gaussian assumptions for class features, enabling training-free integration of visual and text	×	0.03
Ge et al. [220] improve accuracy by identifying ambiguous predictions through prompt and transformation consistency, aug	×	0.03
Pratt et al. [125] enhance open-vocabulary image classification generating discriminative prompts with LLMs, improving a	×	0.05
Menon et al. [126] enhance VLM-based classification by querying LLMs for descriptive features.	×	0.01
Parashar et al. [127] leverage LLMs to identify frequent concept synonyms in pretraining data, enabling more effective p	×	0.02
Shu et al. [128] dynamically adjust prompts for each test sample by minimizing entropy across augmented views.	×	0.05
Feng et al. [130] use pre-trained diffusion models for diverse data augmentation and cosine similarity-based filtration	×	0.06
Ma et al. [132] introduce a self-supervised contrastive learning framework with dual prompts—an online prompt and a hist	×	0.03
Osowiechi et al. [129] improve prediction accuracy by leveraging weight averaging with different text prompts and incorp	×	0.04
Farina et al. [131] improve generalization by augmenting predictions, retaining only the most confident ones, and margin	×	0.02
Rao et al. [135] introduce a framework that repurposes pre-trained CLIP knowledge by transforming image-text matching in	×	0.06
Zhou et al. [133] enhance pixel-level dense prediction by integrating CLIP embeddings with pseudo-labeling and self-trai	×	0.01
Zhang et al. [134] improve the robustness and efficiency of SAM for image segmentation under significant distribution sh	×	0.02

References

- <http://arxiv.org/abs/2501.18592v4>
- <http://arxiv.org/abs/2410.12381v3>
- <http://arxiv.org/abs/2407.04973v1>