

Region-Specific Visual Artifacts and Zero-Shot Dialect Classification in Vision-Language Models

Assignee Research

June 8, 2026

Abstract

This report synthesises findings from 10 peer-reviewed papers addressing the following research question: How does the integration of region-specific visual artifacts impact the zero-shot dialect classification accuracy of vision-language models on the AraDiCE benchmark compared to text-only baselines. 12 claims were extracted from source literature; 3 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Zero-Shot Visual Reasoning by Vision-Language Models: Benchmarking and Analysis. Research question: How does the integration of region-specific visual artifacts impact the zero-shot dialect classification accuracy of vision-language models on the AraDiCE benchmark compared to text-only baselines?.

2 Methodology

Systematic literature search across multiple databases yielded 10 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.5/10.

3 Results

10 papers retrieved. 12 claims extracted; 3 independently verified. Quality review score: 5.5/10.

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
The CLEVR and PTR datasets contain questions requiring minimal world knowledge but a broader range of reasoning steps an	✓	0.21
The CLEVR and PTR datasets provide detailed meta-information for each (question, image) pair, including a complete symbo	×	0.03
LLMs consistently outperform VLMs that utilize the same base LLMs.	×	0.04
Using only its base LLM, i.e. Flan-T5, without the visual front-end achieves $\sim 18\%$ higher accuracy on the PTR dataset.	×	0.10
GPT-4 was $\sim 17\%$ more accurate than GPT-4V on CLEVR.	×	0.02
For questions which can be solved in 2 to 5 reasoning steps, LLMs show performance levels which are significantly above	×	0.09
CoT prompting for visual reasoning in LLMs only obtains better results than standard prompting at large model scales (in	✓	0.20
Using synthetic datasets to benchmark VLMs which are not explicitly trained on reasoning on synthetically rendered scene	×	0.10
LLMs receiving only ground-truth textual scene information consistently perform better than when provided with visual em	✓	0.20
Flan-T5-XL (3B) only performed marginally worse than its larger 11B cousin.	×	0.01
The code uses 2 major libraries for the experiments: the huggingface transformers library for LLM experiments and the Sa	×	0.04
The 2 major datasets used (CLEVR, PTR and GQA) can be downloaded from specific links.	×	0.02

References

- <http://arxiv.org/abs/2404.07214v4>
- <http://arxiv.org/abs/2409.00106v1>
- <http://arxiv.org/abs/2409.11404v3>