

# FineT2I Text-Image Alignment and Zero-Shot CLIP Performance vs. LAION-400M

Assignee Research

June 8, 2026

## Abstract

This report synthesises findings from 11 peer-reviewed papers addressing the following research question: How does the Fine-T2I dataset's text-image alignment quality impact the zero-shot classification accuracy of CLIP models compared to LAION-400M. 8 claims were extracted from source literature; 8 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 8.9/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: MobileCLIP: Fast Image-Text Models through Multi-Modal Reinforced Training. Research question: How does the Fine-T2I dataset's text-image alignment quality impact the zero-shot classification accuracy of CLIP models compared to LAION-400M?.

## 2 Methodology

Systematic literature search across multiple databases yielded 11 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 8.9/10.

## 3 Results

11 papers retrieved. 8 claims extracted; 8 independently verified. Quality review score: 8.9/10.

## 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
Contrastive pretraining of image-text foundation models, such as CLIP, demonstrated excellent zero-shot performance and	✓	0.34
These models utilize large transformer-based encoders with significant memory and latency overhead which pose challenges	✓	0.32
MobileCLIP is a new family of efficient image-text models optimized for runtime performance along with a novel and effic	✓	0.43
The proposed training approach leverages knowledge transfer from an image captioning model and an ensemble of strong CLI	✓	0.37
MobileCLIP sets a new state-of-the-art latency-accuracy tradeoff for zero-shot classification and retrieval tasks on sev	✓	0.31
MobileCLIP-S2 variant is 2.3 $\times$ faster while more accurate compared to previous best CLIP model based on ViT-B/16.	✓	0.36
Multi-modal reinforced training achieves +2.9% average performance improvement on 38 evaluation benchmarks compared to t	✓	0.35
The proposed approach achieves 10 $\times$ -1000 $\times$ improved learning efficiency when compared with non-reinforced CLIP training.	✓	0.29

## References

- <https://doi.org/10.48550/arxiv.2307.12980>
- <https://doi.org/10.48550/arxiv.2206.02770>
- <https://doi.org/10.48550/arxiv.2311.17049>