

Scaling Soft-Labeled Synthetic Data for Zero-Shot LLM Reasoning in CALVIN

Assignee Research

June 7, 2026

Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: To what extent does scaling the size of soft-labeled synthetic datasets affect the zero-shot reasoning capabilities of LLMs in the CALVIN task, quantified by improvements in success rate and human. 0 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: An Evaluation of Large Pre-Trained Models for Gesture Recognition using Synthetic Videos. Research question: To what extent does scaling the size of soft-labeled synthetic datasets affect the zero-shot reasoning capabilities of LLMs in the CALVIN task, quantified by improvements in success rate and human preference alignment metrics?.

2 Methodology

Systematic literature search across multiple databases yielded 16 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.3/10.

3 Results

16 papers retrieved. 0 claims extracted; 0 independently verified. Quality review score: 4.3/10.

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.

References

- <http://arxiv.org/abs/2410.02152v1>
- <http://arxiv.org/abs/2304.10464v4>
- <http://arxiv.org/abs/2403.09832v1>