

Scaling Unlabeled Video-Audio Pretraining for Few-Shot Latent Action Model Adaptation

Assignee Research

June 8, 2026

Abstract

This report synthesises findings from 10 peer-reviewed papers addressing the following research question: How does increasing the scale of unlabeled video-audio pretraining data impact the few-shot adaptation accuracy of latent action models on the RoboBench benchmark compared to supervised baselines. 12 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: CLAM: Continuous Latent Action Models for Robot Learning from Unlabeled Demonstrations. Research question: How does increasing the scale of unlabeled video-audio pretraining data impact the few-shot adaptation accuracy of latent action models on the RoboBench benchmark compared to supervised 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 4.2/10.

3 Results

10 papers retrieved. 12 claims extracted; 1 independently verified. Quality review score: 4.2/10.

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
CLAM outperforms all baselines and nearly matches the performance of BC with expert data in both state- and image-based	×	0.05
CLAM improves upon the best baseline VPT by more than 2 \times average normalized return on the DMControl (locomotion) tasks.	×	0.07
CLAM improves upon the best baseline VPT by around 2-3 \times success rate on the MetaWorld (manipulation) tasks.	×	0.11
Transformer-CLAM achieves performance close to or even better than that of BC-Expert which uses the same amount of privi	×	0.08
All variants of CLAM outperform the best baseline VPT.	×	0.05
CLAM outperforms state-of-the-art methods in the problem setting where only play data is available as action-labeled dat	✓	0.17
CLAM scales with Dunlabeled while supervised IDMs only scale with Dlabeled .	×	0.03
CLAM can leverage vast, unstructured observation data to learn latent actions in an unsupervised manner.	×	0.11
CLAM enables scalable learning from easy-to-collect, cheap play data avoiding the need for expensive task-specific data	×	0.04
The Transformer-CLAM model uses 6 encoder layers, 6 decoder layers, a feedforward dimension of 2048, 8 attention heads,	×	0.01
The CALVIN environment has a maximum of 200 episode steps, a state dimension of 39, an action dimension of 7, an image s	×	0.03
The MetaWorld environment has a maximum of 100 episode steps, a state dimension of 39, an action dimension of 4, an imag	×	0.04

References

- <http://arxiv.org/abs/2402.04177v3>
- <http://arxiv.org/abs/2505.04999v1>

- <http://arxiv.org/abs/1911.06045v3>