

# What is the comparative robustness of continuous latent action models versus discrete token methods when learn

Assignee Research

June 10, 2026

## Abstract

Learning robot policies using imitation learning requires collecting large amounts of costly action-labeled expert demonstrations, which fundamentally limits the scale of training data. A promising approach to address this bottleneck is to harness the abundance of unlabeled observations-e.g., from video demonstrations-to learn latent action labels in an unsupervised way. However, we find that existing methods struggle when applied to complex robot tasks requiring fine-grained motions. We design continuous latent action models (CLAM) which incorporate two key ingredients we find necessary for l

## 1 Introduction

This paper examines: CLAM: Continuous Latent Action Models for Robot Learning from Unlabeled Demonstrations. Research question: What is the comparative robustness of continuous latent action models versus discrete token methods when learning from unlabeled video demonstrations with varying levels of observation noise and occlusion?.

## 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 4.0/10.

## 3 Results

10 papers retrieved. 19 claims extracted; 0 independently verified. Quality review score: 4.0/10.

## 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
CLAM improves upon the best baseline VPT by more than 2 $\times$ average normalized return on DMControl locomotion tasks.	×	0.08
CLAM improves upon the best baseline VPT by around 2-3 $\times$ success rate on MetaWorld manipulation tasks.	×	0.11
On the HalfCheetah task, TF-CLAM achieved a normalized return of 0.72 $\pm$ 0.04.	×	0.02
On the HalfCheetah task, BC-Expert achieved a normalized return of 0.68 $\pm$ 0.02.	×	0.02
On the Hopper task, TF-CLAM achieved a normalized return of 0.81 $\pm$ 0.05.	×	0.02
On the Hopper task, BC-Expert achieved a normalized return of 0.76 $\pm$ 0.04.	×	0.02
Transformer-CLAM achieves performance close to or better than BC-Expert in several tasks.	×	0.04
The Transformer CLAM model uses 6 encoder layers and 6 decoder layers.	×	0.04
The Transformer CLAM model uses a feedforward dimension of 2048 and 4 attention heads.	×	0.01
The CALVIN Transformer CLAM model uses 8 attention heads.	×	0.02
The MetaWorld environment configuration uses a state dimension of 39 and an action dimension of 4.	×	0.04
The MetaWorld environment configuration uses an image shape of [84, 84, 3] and stacks 3 frames.	×	0.01
The CALVIN environment configuration uses a state dimension of 39 and an action dimension of 7.	×	0.04
The CALVIN environment configuration uses an action repeat of 7.	×	0.03
The study evaluates locomotion tasks (Hopper and HalfCheetah) from the DMControl benchmark.	×	0.03
The study evaluates manipulation tasks (Assembly, Bin Picking, Peg Insert Side, and Shelf Place) from the MetaWorld benchmark.	×	0.04
The study evaluates Close Drawer and Slider Left tasks in the CALVIN environment.	×	0.01
All domains used in the experiments are continuous control environments with fixed episode lengths and no termination cost.	×	0.06
BC-AL performs poorly because it imitates sub-optimal demonstrations.	×	0.03

## References

- <http://arxiv.org/abs/2302.08893v4>
- <http://arxiv.org/abs/2505.04999v1>
- <http://arxiv.org/abs/2605.15725v1>