

Scaling Motion-Aware Fine-Tuning Data for MoCLIP Robustness in Adversarial Text-to-Motion Generation

Assignee Research

June 16, 2026

Abstract

Human motion generation is essential for fields such as animation, robotics, and virtual reality, requiring models that effectively capture motion dynamics from text descriptions. Existing approaches often rely on Contrastive Language-Image Pretraining (CLIP)-based text encoders, but their training on text-image pairs constrains their ability to understand temporal and kinematic structures inherent in motion and motion generation. This work introduces MoCLIP, a fine-tuned CLIP model with an additional motion encoding head, trained on motion sequences using contrastive learning and tethering lo

1 Introduction

This paper examines: MoCLIP: Motion-Aware Fine-Tuning and Distillation of CLIP for Human Motion Generation. Research question: What is the impact of scaling the motion-aware fine-tuning dataset size on MoCLIP’s robustness to adversarial perturbations in text-to-motion generation, as measured by FID scores on Action3D?.

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 8.7/10.

3 Results

11 papers retrieved. 11 claims extracted; 10 independently verified. Quality review score: 8.7/10.

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
MoCLIP improves Top-1, Top-2, and Top-3 accuracy while maintaining competitive FID, leading to improved text-to-motion a	✓	0.34
The motion encoder used in MoCLIP generates robust motion embeddings with strong semantic coherence.	✓	0.16
MoCLIP introduces cross-limb attention connections that extend beyond conventional skeletal adjacency constraints.	✓	0.17
MoCLIP includes direct attention connections between both hands and both feet to better capture inter-limb coordination	✓	0.23
Temporal attention mechanisms are applied to the encoded motion features before pooling along the temporal dimension in	✓	0.19
MoCLIP uses a multi-term loss function to achieve effective contrastive alignment, preserve original semantic representa	✓	0.25
MoCLIP employs a symmetric cross-entropy loss following standard contrastive learning practice.	✓	0.17
MoCLIP uses a feature distillation loss inspired by recent works in CLIP fine-tuning, such as CLIP-CITE and LDIFS.	✓	0.22
The proposed model relies on pre-trained weights from each chosen baseline model on HumanML3D and KIT-ML datasets.	✓	0.20
MoCLIP proposes a specialized fine-tuning strategy for the CLIP graph-based human motion data.	✓	0.21
MoCLIP fine-tunes the textual embeddings using a distillation loss.	×	0.15

References

- <http://arxiv.org/abs/2602.09439v1>
- <http://arxiv.org/abs/2505.10810v1>
- <http://arxiv.org/abs/2205.14697v1>