

# Performance of Large Pre-trained Gesture Recognition Models on Real-world Datasets after Synthetic Fine-tuning

Assignee Research

June 13, 2026

## Abstract

In this work, we explore the possibility of using synthetically generated data for video-based gesture recognition with large pre-trained models. We consider whether these models have sufficiently robust and expressive representation spaces to enable "training-free" classification. Specifically, we utilize various state-of-the-art video encoders to extract features for use in k-nearest neighbors classification, where the training data points are derived from synthetic videos only. We compare these results with another training-free approach – zero-shot classification using text descriptions o

## 1 Introduction

This paper examines: An Evaluation of Large Pre-Trained Models for Gesture Recognition using Synthetic Videos. Research question: How do large pre-trained models for gesture recognition perform when evaluated on real-world video datasets after fine-tuning on synthetically generated videos, measured by classification accuracy and generalization to unseen gestures?.

## 2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.6/10.

## 3 Results

13 papers retrieved. 18 claims extracted; 14 independently verified. Quality review score: 7.6/10.

## 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
The RoCoG-v2 dataset consists of 7 gesture categories.	✓	0.16
The synthetic training data consists of 44K videos.	×	0.13
The small real dataset available for training consists of 203 videos.	×	0.13
The K value used for KNN classification is 3.	×	0.08
The UMT model is pre-trained on K710 videos.	✓	0.19
The ViCLIP model is pre-trained on a filtered version of the InternVid dataset with 10M video-text pairs.	✓	0.19
The VideoMAE models are pre-trained on a larger dataset of 1.3B videos.	×	0.15
The KNN accuracy for ViT-B/16 with UMT pre-training on K710 data is 18.2% for synthetic train and 31.2% for real train.	✓	0.18
The KNN accuracy for ViT-B/16 with ViCLIP pre-training on InternVid FLT-10M data is 19.2% for synthetic train and 40.4%	✓	0.20
The KNN accuracy for ViT-B/16 with UMT pre-training on K710 data and fine-tuning on K710 data is 42.4% for synthetic tra	✓	0.19
The KNN accuracy for ViT-B/16 with UMT pre-training on K710 data and fine-tuning on K710 + K400 data is 38.4% for synthe	✓	0.20
The KNN accuracy for ViT-B/16 with UMT pre-training on K710 data and fine-tuning on K710 + K600 data is 33.3% for synthe	✓	0.21
The KNN accuracy for ViT-B/16 with UMT pre-training on K710 data and fine-tuning on K710 + K700 data is 35.4% for synthe	✓	0.20
The KNN accuracy for ViT-B/16 with VideoMAE pre-training on UnlabeledHybrid data and fine-tuning on K710 data is 32.3% f	✓	0.19
The KNN accuracy for ViT-B/16 with VideoMAE pre-training on UnlabeledHybrid data and fine-tuning on SSv2 data is 43.4% f	✓	0.18
The KNN accuracy for ViT-L/16 with VideoMAE pre-training on UnlabeledHybrid data and fine-tuning on SSv2 data is 64.6% f	✓	0.18
The zero-shot classification accuracy for ViCLIP-B with InternVid FLT-10M training data and original text descriptions i4	✓	0.27
The zero-shot classification accuracy for ViCLIP-B with InternVid FLT-10M training data and transformed text description	✓	0.25

## References

- <http://arxiv.org/abs/2410.02152v1>
- <http://arxiv.org/abs/2404.17929v1>
- <http://arxiv.org/abs/2412.01508v1>