

Video-JEPA Performance with Dynamic vs. Fixed Auxiliary Loss Weighting on Kinetics-400

Assignee Research

June 8, 2026

Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: How does the performance of Video-JEPA with dynamic auxiliary loss weighting compare to fixed weighting schemes on the Kinetics-400 benchmark under identical fine-tuning conditions. 14 claims were extracted from source literature; 4 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Deep Analysis of CNN-based Spatio-temporal Representations for Action Recognition. Research question: How does the performance of Video-JEPA with dynamic auxiliary loss weighting compare to fixed weighting schemes on the Kinetics-400 benchmark under identical fine-tuning conditions?.

2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 5.2/10.

3 Results

12 papers retrieved. 14 claims extracted; 4 independently verified. Quality review score: 5.2/10.

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
Kinetics and Something-Something are popular action recognition datasets.	×	0.09
CNN-based approaches have made impressive progress in action recognition.	✓	0.17
There are several fundamental questions that still largely remain unanswered in the field of action recognition.	×	0.05
Kinetics-400 has recently emerged as a primary benchmark for action recognition.	×	0.10
Kinetics-400 is known to be strongly biased towards spatial modeling.	×	0.05
I3D based on 3D-InceptionV1 has become a 'gatekeeper' baseline to compare with for any recently proposed approaches of a	×	0.09
I3D, with ResNet50 as backbone, performs comparably with or outperforms many recent methods that are claimed to be better	×	0.04
Performance evaluation in action recognition may be further confounded by many other issues such as variations in training	×	0.05
Temporal pooling tends to suppress the efficacy of temporal modeling in an action model.	×	0.04
Temporal pooling provides a significant performance boost to TSN.	×	0.05
By removing non-structural differences between 2D-CNN and 3D-CNN models, they behave similarly in terms of spatio-temporal	✓	0.34
The paper presents a unified framework for 2D-CNN and 3D-CNN approaches.	✓	0.25
The paper systematically compares 2D-CNN and 3D-CNN models to better understand the differences and spatio-temporal behavior	✓	0.26
The paper thoroughly benchmarks several SOTA approaches and compares them with I3D.	×	0.03

References

- <http://arxiv.org/abs/2605.17165v1>
- <http://arxiv.org/abs/2010.11757v4>

- <http://arxiv.org/abs/2604.12406v2>