

# Video-JEPA Auxiliary Objectives and Frozen Representation Accuracy in Video Benchmarks

Assignee Research

June 7, 2026

## Abstract

This report synthesises findings from 10 peer-reviewed papers addressing the following research question: How do different auxiliary objective variants in Video-JEPA affect frozen representation accuracy on downstream video benchmarks. 16 claims were extracted from source literature; 3 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Factorized Latent Dynamics for Video JEPA: An Empirical Study of Auxiliary Objectives. Research question: How do different auxiliary objective variants in Video-JEPA affect frozen representation accuracy on downstream video benchmarks?.

## 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 4.3/10.

## 3 Results

10 papers retrieved. 16 claims extracted; 3 independently verified. Quality review score: 4.3/10.

## 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
Motion-Guided Masking improves Diving-48 accuracy by 0.30 percentage points in the UCF-101 pretraining setting.	×	0.14
Motion-Guided Masking improves ImageNet-100 accuracy by 0.14 percentage points in the UCF-101 pretraining setting.	×	0.14
Motion-Guided Masking improves SSv2 accuracy by 1.38 percentage points in the UCF-101 pretraining setting.	×	0.09
Kinematic variants degrade Diving-48 accuracy by 2.5 to 2.9 points in the UCF-101 pretraining setting.	×	0.10
Kinematic variants improve ImageNet-100 accuracy by 1.5 to 1.7 points in the UCF-101 pretraining setting.	×	0.11
FWM-HW-LD achieves a +5.92 percentage point gain on ImageNet-100 in mixed-dataset pretraining.	✓	0.16
FWM-HW-LD achieves a +3.21 percentage point gain on SSv2 in mixed-dataset pretraining.	✓	0.15
FWM-HW-LD results in a -0.30 percentage point change on Diving-48 relative to the reference baseline in mixed-dataset pr	✓	0.20
LD-JEPA achieves a +5.02 percentage point gain on SSv2 in mixed-dataset pretraining.	×	0.09
10 out of 14 tested methods lose more than 5 points on ImageNet-100 in the mixed-dataset setup.	×	0.08
AC-JEPA and FAC-JEPA objectives result in performance losses of 13 to 16 percentage points on ImageNet-100.	×	0.10
In the ablation study, LD alone boosts SSv2 accuracy by +5.02 points but hurts ImageNet and Diving-48.	×	0.08
In the ablation study, FWM alone boosts ImageNet accuracy by +1.88 points but hurts SSv2 and Diving-48.	×	0.09
The combination of FWM+LD without hard weighting performs poorly on ImageNet with a -10.14 point change.	×	0.06
Synthetic Motion Discrimination shows a +40 to +45 point improvement with kinematic regularization.	×	0.01
The encoder produces a fixed 768-dimensional embedding.	×	0.03

## References

- <http://arxiv.org/abs/2605.17165v1>
- <http://arxiv.org/abs/2304.00571v3>
- <http://arxiv.org/abs/2404.08471v1>