

# OpenPangu-7B-MLA Robustness to Noisy Video Inputs in Multimodal Action Recognition

Assignee Research

June 9, 2026

## Abstract

This report synthesises findings from 4 peer-reviewed papers addressing the following research question: How does the robustness of OpenPangu-7B-MLA to noisy video inputs on Mini Kinetics-C compare to other multimodal models in terms of action recognition accuracy. 16 claims were extracted from source literature; 6 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.0/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Advancing multimodal emotion recognition in big data through prompt engineering and deep adaptive learning. Research question: How does the robustness of OpenPangu-7B-MLA to noisy video inputs on Mini Kinetics-C compare to other multimodal models in terms of action recognition accuracy?.

## 2 Methodology

Systematic literature search across multiple databases yielded 4 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 6.0/10.

## 3 Results

4 papers retrieved. 16 claims extracted; 6 independently verified. Quality review score: 6.0/10.

## 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
The framework integrates text, audio, video, and motion data.	✓	0.20
The framework employs Generative Adversarial Networks (GANs) for synthetic sample generation to address class imbalance.	✓	0.24
The framework uses Dynamic Prompt Engineering (DPE) for enhanced feature extraction across modalities.	✓	0.22
Text features are processed using the Mistral-7B model.	×	0.14
Audio features are processed using the HuBERT model.	×	0.10
Video features are processed using TimeSformer and LLaVA.	×	0.12
Motion features are processed using MediaPipe Pose.	×	0.12
The system fuses inputs using Hierarchical Attention-based Graph Neural Networks (HAN-GNN).	✓	0.24
The system fuses inputs using Cross-Modality Transformer Fusion (XMTF).	✓	0.18
The framework utilizes contrastive learning with Prototypical Networks to enhance class separation.	✓	0.19
The framework achieved a training accuracy of 99.92% on the IEMOCAP dataset.	×	0.10
The framework achieved a training accuracy of 99.95% on the MELD dataset.	×	0.10
The framework achieved a testing accuracy of 99.82% on the IEMOCAP dataset.	×	0.09
The framework achieved a testing accuracy of 99.81% on the MELD dataset.	×	0.09
Training for the framework was completed in 5 minutes.	×	0.05
Inference time for the framework is under 0.4 ms per sample.	×	0.10

## References

- <https://doi.org/10.48550/arxiv.2304.00685>

- <https://doi.org/10.1186/s40537-025-01264-w>
- <https://doi.org/10.48550/arxiv.2305.06324>