

# Causal Encoder Integration in WALT Enhancing Multimodal Video Captioning Performance

Assignee Research

May 30, 2026

## Abstract

This report synthesises findings from 11 peer-reviewed papers addressing the following research question: How does the causal encoder design in W.A.L.T impact downstream performance when integrated with state-of-the-art multimodal models like Flamingo or PaLI on video captioning benchmarks (e.g., Multimodal learning on video and text has seen significant progress, particularly in tasks like text-to-video retrieval, video-to-text retrieval, and video captioning. However, most existing methods and datasets focus exclusively on English. 18 claims were extracted from source literature; 5 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 5.6/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: MSVD-Indonesian: A Benchmark for Multimodal Video-Text Tasks in Indonesian. Research question: How does the causal encoder design in W.A.L.T impact downstream performance when integrated with state-of-the-art multimodal models like Flamingo or PaLI on video captioning benchmarks (e.g., ActivityNet, MSRVTT)?.

## 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 5.6/10.

## 3 Results

11 papers retrieved. 18 claims extracted; 5 independently verified. Quality review score: 5.6/10.

## 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 MSVD-Indonesian dataset improves performance across all tasks and metrics.	✓	0.21
The MSVD-Indonesian dataset is released publicly to support future research.	✓	0.19
The standard split of the MSVD dataset is 1200, 100, and 670 videos for train, validation, and test set.	×	0.06
Evaluation metrics for retrieval include R@1, R@5, R@10, MedianRank, and MeanRank.	×	0.05
Evaluation metrics for captioning include BLEU@4, ROUGE-L, METEOR, and CIDEr.	×	0.04
The feature extractor for X-CLIP is a pretrained CLIP (ViT-B/16) model.	×	0.03
The learning rate for X-CLIP is set to 1e-4.	×	0.03
Hyperparameters for X-CLIP include maximum word length of 32, maximum frame length of 12, and 5 training epochs.	×	0.01
The batch size for X-CLIP training is 16.	×	0.03
Training X-CLIP takes about 15 hours on a Linux environment with 1 NVIDIA GeForce GTX 1080 Ti.	×	0.01
Video features for VNS-GRU are extracted using Efficient Convolutional Network (ECN) pretrained on Kinetics-400 dataset.	×	0.05
The MSVD-Indonesian dataset is the first Indonesian video-text dataset, translated from MSVD.	✓	0.24
Baseline results are established for three tasks using models originally developed for English video-text tasks.	✓	0.18
Cross-lingual transfer learning is effective for Indonesian video-text tasks.	✓	0.25
Most existing video-text datasets, including MSVD, MSR-VTT, and ActivityNet Captions, are constructed with English annot	×	0.13
A few multilingual datasets exist, including MSVD-CN, MSVD-Turkish, and others in Chinese, Hindi, and Italian.	×	0.03
The original MSVD dataset was collected in multiple languages, but the publicly released version includes only English.	×	0.07
No Indonesian version of the MSVD dataset existed before the release of MSVD-Indonesian.	×	0.12

## References

- <http://arxiv.org/abs/1806.08854v1>
- <http://arxiv.org/abs/2306.11341v2>
- <http://arxiv.org/abs/1911.12018v6>