

Hierarchical Video-Subtitle Matching in HERO for Cross-Lingual Video Retrieval

Assignee Research

June 8, 2026

Abstract

This report synthesises findings from 11 peer-reviewed papers addressing the following research question: How does the Video-Subtitle Matching objective in HERO perform on cross-lingual video retrieval tasks compared to standard contrastive learning baselines. 11 claims were extracted from source literature; 4 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: HERO: Hierarchical Encoder for Video+Language Omni-representation Pre-training. Research question: How does the Video-Subtitle Matching objective in HERO perform on cross-lingual video retrieval tasks compared to standard contrastive learning baselines?.

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 6.2/10.

3 Results

11 papers retrieved. 11 claims extracted; 4 independently verified. Quality review score: 6.2/10.

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
HERO achieves new state of the art on multiple benchmarks including Text-based Video/Video-moment Retrieval, Video Quest	✓	0.39
HERO introduces two new challenging benchmarks: How2QA for Video QA and How2R for Video Retrieval, collected from divers	✓	0.26
HERO is evaluated on 6 existing benchmarks: TVR, TVQA, VIOLIN, TVC, DiDeMo, and MSR-VTT.	×	0.03
HERO’s pre-training dataset is composed of 7.6M video clips with their accompanying subtitles from TV and HowTo100M data	×	0.08
The TV Dataset contains 21,793 video clips from 925 episodes, each 60-90 seconds long, covering long-range scenes with c	×	0.05
The HowTo100M Dataset contains 1.22 million videos, with activities falling into 12 categories, and each video is associ	×	0.02
The pre-processed subset of HowTo100M consists of 7.56M video clips, each 60 seconds long, accompanied with English subt	×	0.03
How2R is a new benchmark for text-based video-moment retrieval, collected from HowTo100M videos using Amazon Mechanical	×	0.12
How2QA is a new benchmark for video question answering, collected from HowTo100M videos.	×	0.11
HERO encodes multimodal inputs in a hierarchical structure, capturing local context via Cross-modal Transformer and glob	✓	0.30
HERO is jointly trained on HowTo100M and large-scale TV datasets to gain deep understanding of complex social dynamics w	✓	0.33

References

- <http://arxiv.org/abs/2005.00200v2>
- <http://arxiv.org/abs/2506.15415v1>
- <http://arxiv.org/abs/2008.02531v2>