

# Domain-Specific Videoqa Models Perform On Out-Of-Distribution Tasks In The Howto100M Benchmark When Compared To Models

Assignee Research

June 8, 2026

## Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: How do domain-specific VideoQA models perform on out-of-distribution tasks in the HowTo100M benchmark when compared to models trained with generalized pre-training objectives. 18 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Revisiting Out-of-distribution Robustness in NLP: Benchmark, Analysis, and LLMs Evaluations. Research question: How do domain-specific VideoQA models perform on out-of-distribution tasks in the HowTo100M benchmark when compared to models trained with generalized pre-training objectives?.

## 2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.7/10.

## 3 Results

14 papers retrieved. 18 claims extracted; 0 independently verified. Quality review score: 3.7/10.

## 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
Previous work on OOD robustness in NLP lacks standardized benchmark suites, leading to heuristic and popularity-based da	×	0.10
The proposed protocol requires ID datasets to be large and diverse for comprehensive knowledge.	×	0.05
The proposed protocol prioritizes OOD datasets with distinct distributions and dissimilarity regarding text sources and	×	0.03
The proposed protocol prioritizes challenging distribution shifts based on performance degradation.	×	0.11
The BOSS benchmark suite covers sentiment analysis, toxic detection, natural language inference, name entity recognition	×	0.09
The BOSS benchmark establishes one ID dataset and three corresponding OOD datasets for each task type.	×	0.09
The Yelp Product dataset contains 30,000 training samples and 38,905 test samples with an average training length of 71.	×	0.04
The Dynasent dataset is categorized as an Adversarial source with 93,553 training samples and 4,320 test samples.	×	0.05
The semantic similarity score between the Amazon and Yelp datasets is 49.22.	×	0.03
The semantic similarity score between the IMDB and SST datasets is 84.62.	×	0.01
For the Sentiment Analysis (SA) task, the ID dataset is Amazon (AZ).	×	0.03
For the Sentiment Analysis (SA) task, the OOD datasets are Dynasent (DS), SemEval (SE), and SST (SST).	×	0.03
For the Toxic Detection (TD) task, the ID dataset is Civil Comments (CC).	×	0.03
For the Toxic Detection (TD) task, the OOD datasets are AdvCivil (AC), Implicit Hate (IH), and ToxiGen (TG).	×	0.02
For the Natural Language Inference (NLI) task, the ID dataset is MNLI (MN).	×	0.05
For the Named Entity Recognition (NER) task, the ID dataset is FewNerd (FN).	×	0.02
For the Extractive Question Answering (EQA) task, the ID dataset is SQuAD. 4	×	0.09
The performance score on the SQuAD ID dataset is 38, while the score on the AdvQA OOD dataset is lower.	×	0.04

## References

- <http://arxiv.org/abs/2306.04618v2>
- <http://arxiv.org/abs/1905.03197v3>
- <http://arxiv.org/abs/2311.11096v1>