

# Comparison of Multimodal Semantic Similarity Prediction Errors Under Synthetic and Natural Distribution Shifts

Assignee Research

June 11, 2026

## Abstract

Deep learning (DL) models are widely used in real-world applications but remain vulnerable to distribution shifts, especially due to weather and lighting changes. Collecting diverse real-world data for testing the robustness of DL models is resource-intensive, making synthetic corruptions an attractive alternative for robustness testing. However, are synthetic corruptions a reliable proxy for real-world corruptions? To answer this, we conduct the largest benchmarking study on semantic segmentation models, comparing performance on real-world corruptions and synthetic corruptions datasets. Our r

## 1 Introduction

This paper examines: Are Synthetic Corruptions A Reliable Proxy For Real-World Corruptions?. Research question: How does the mean squared error of multimodal semantic similarity predictions on the SICK-R entailment subset compare between synthetic image corruptions and natural distribution shifts in lighting and weather conditions?.

## 2 Methodology

Systematic literature search across multiple databases yielded 8 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.2/10.

## 3 Results

8 papers retrieved. 10 claims extracted; 7 independently verified. Quality review score: 7.2/10.

## 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
ACDC (Adverse Conditions Dataset with Correspondences) consists of images from similar regions and scenes as Cityscapes	✓	0.27
There is a very strong positive correlation between performance against synthetic corruptions from 2D Common Corruptions	✓	0.21
The Pearson correlation between mean performance on 2D Common Corruptions and ACDC mIoU is 0.795.	✓	0.19
The Pearson correlation between Cityscapes GAM3 (worst-case scenario) and ACDC mIoU is 0.828.	×	0.13
The Pearson correlation between mean performance against all 2D Common Corruptions and performance against worst-case co	✓	0.25
Synthetic Snow corruption shows a Pearson correlation of 0.867 with real-world snow-related degradation in the ACDC data	✓	0.21
Synthetic Brightness corruption exhibits a weak alignment with real-world conditions, showing a Pearson correlation of 0	×	0.15
Synthetic Fog corruption exhibits a weak alignment with real-world atmospheric distortions, showing a Pearson correlatio	×	0.12
Common Corruptions and 3D Common Corruptions are tools proposed for benchmarking the robustness of image classification	✓	0.23
The ACDC dataset serves as a community-accepted tool for benchmarking real-world OOD robustness.	✓	0.22

## References

- <http://arxiv.org/abs/2107.12052v2>
- <http://arxiv.org/abs/2112.03057v1>
- <http://arxiv.org/abs/2505.04835v1>