

# ViLT Model Robustness in Adversarial VQA with SMOTE and StyleGAN Augmentation

Assignee Research

June 12, 2026

## Abstract

Machine learning (ML) systems have introduced significant advances in various fields, due to the introduction of highly complex models. Despite their success, it has been shown multiple times that machine learning models are prone to imperceptible perturbations that can severely degrade their accuracy. So far, existing studies have primarily focused on models where supervision across all classes were available. In contrast, Zero-shot Learning (ZSL) and Generalized Zero-shot Learning (GZSL) tasks inherently lack supervision across all classes. In this paper, we present a study aimed on evaluat

## 1 Introduction

This paper examines: A Deep Dive into Adversarial Robustness in Zero-Shot Learning. Research question: How does the robustness of ViLT models on the Adversarial VQA benchmark compare when trained with SMOTE versus StyleGAN-based data augmentation?.

## 2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 7.7/10.

## 3 Results

12 papers retrieved. 10 claims extracted; 8 independently verified. Quality review score: 7.7/10.

## 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
The CUB dataset contains 312 attributes, 200 classes, and 11788 images.	✓	0.16
The SUN dataset contains 102 attributes, 717 classes, and 14340 images.	✓	0.15
The AWA2 dataset contains 85 attributes, 50 classes, and 37322 images.	×	0.15
The standard per-class top-1 accuracy is used for ZSL evaluation.	✓	0.16
For GZSL, per-class top-1 accuracy values for seen and unseen classes are used to compute harmonic-scores.	✓	0.23
The reproduced values of ALE are denoted as original, although there are slight variations compared to the original resu	✓	0.15
The ALE model is formulated as $F(x, y; W) = \theta(x)W^T \varphi(y)$ , where $\theta(x)$ is the visual and $\varphi(y)$ is the class embeddings.	✓	0.19
The ALE model is one of the earlier studies that showed direct mapping by exploiting data and auxiliary information is m	✓	0.28
The ALE model is a stable and competitive model in modern benchmarks [57].	✓	0.20
The ALE model is selected for evaluating adversarial robustness in ZSL approaches.	×	0.14

## References

- <http://arxiv.org/abs/2104.09630v2>
- <http://arxiv.org/abs/2110.06166v4>
- <http://arxiv.org/abs/2008.07651v1>