

# Diffusion Models vs. VAEs in Anomaly Detection: Precision-Recall Performance on OpenML Benchmarks

Assignee Research

June 9, 2026

## Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: How does the precision-recall performance of diffusion models compare to variational autoencoders (VAEs) on OpenML anomaly detection benchmarks when evaluated under varying contamination rates. 17 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Robust outlier detection by de-biasing VAE likelihoods. Research question: How does the precision-recall performance of diffusion models compare to variational autoencoders (VAEs) on OpenML anomaly detection benchmarks when evaluated under varying contamination rates?.

## 2 Methodology

Systematic literature search across multiple databases yielded 16 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.7/10.

## 3 Results

16 papers retrieved. 17 claims extracted; 1 independently verified. Quality review score: 4.7/10.

## 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 study trained and tested VAEs on four grayscale image datasets: MNIST, Fashion-MNIST, EMNIST, and Sign Language-MNIS	×	0.05
For each grayscale VAE, testing was performed against three other grayscale datasets and one grayscale uniform noise dat	×	0.03
The study trained and tested VAEs on five natural image datasets: SVHN, CelebA, ComprehensiveCars, GTSRB, and CIFAR-10.	×	0.06
For each natural image VAE, testing was performed against four other natural image datasets and one colored uniform nois	×	0.07
Each VAE was trained 6 times using 3 different random initializations across 2 train-validation splits.	×	0.02
Reported performance measures represent the average across all 6 training runs.	×	0.02
Both train and test images were contrast stretched prior to computing bias-corrected likelihood (BC-LL) scores.	×	0.04
Outlier detection performance was evaluated using Area Under the ROC Curve (AUROC).	×	0.07
In nearly every case with grayscale datasets, bias-corrected likelihoods (BC-LL) outperformed uncorrected likelihoods (L	×	0.09
Bias correction enabled the Fashion-MNIST VAE to assign higher likelihoods to in-distribution Fashion-MNIST samples than	×	0.05
The Fashion-MNIST VAE with bias correction achieved a perfect AUROC when distinguishing Fashion-MNIST from MNIST.	×	0.04
Samples with the highest bias-corrected likelihoods were visually more typical than those with the highest uncorrected l	×	0.05
BC-LL scores approached or exceeded state-of-the-art accuracies for outlier detection on multiple grayscale and natural	✓	0.16
VAEs trained on the CIFAR-10 dataset yielded AUROC values ranging from 37 to 66 with bias correction across outlier data	×	0.05
Other approaches based on VAE likelihoods also performed poorly on the CIFAR-10 dataset.	×	0.08
The CIFAR-10 dataset comprises 40 different categories of images that are visually and semantically unrelated.	×	0.02
Four category-specific VAEs were trained on CIFAR-10 categories: Airplane, Ship, Frog, and Deer.	×	0.02

## References

- <http://arxiv.org/abs/1804.06364v2>
- <http://arxiv.org/abs/2108.08760v3>
- <http://arxiv.org/abs/2406.16308v1>