

Scaling Contrastive Learning with Masked Autoencoders in Vision Transformer Point Cloud Generalization

Assignee Research

June 9, 2026

Abstract

This report synthesises findings from 9 peer-reviewed papers addressing the following research question: What is the impact of scaling contrastive learning with masked autoencoder pretraining on the cross-domain generalization capabilities of Vision Transformers (ViTs) in point cloud representation. 14 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Bringing Masked Autoencoders Explicit Contrastive Properties for Point Cloud Self-Supervised Learning. Research question: What is the impact of scaling contrastive learning with masked autoencoder pretraining on the cross-domain generalization capabilities of Vision Transformers (ViTs) in point cloud representation learning?.

2 Methodology

Systematic literature search across multiple databases yielded 9 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.5/10.

3 Results

9 papers retrieved. 14 claims extracted; 2 independently verified. Quality review score: 4.5/10.

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 source code and trained models for Point-CMAE are available at https://github.com/Amazingren/Point-CMAE .	×	0.11
ShapeNet is a synthetic 3D dataset containing 52,470 3D shapes across 55 object categories.	×	0.04
The ShapeNet training set used for pre-training contains 41,952 shapes.	×	0.07
For each 3D shape, 1024 points are sampled to serve as the model input.	×	0.07
Each point cloud is divided into 64 patches (n=64).	×	0.07
The KNN algorithm selects k=32 nearest points to form a point patch.	×	0.01
The proposed method is pre-trained for 300 epochs using an AdamW optimizer.	×	0.05
The encoder in the autoencoder’s backbone consists of 12 Transformer blocks.	×	0.03
The decoder in the autoencoder’s backbone consists of 4 ViT encoder blocks.	×	0.05
Point-CMAE achieves an Overall Accuracy (OA) of 90.02% on ScanObjectNN without rotation data augmentation.	×	0.05
Point-CMAE achieves an Overall Accuracy (OA) of 93.46% on ScanObjectNN with rotation data augmentation.	×	0.05
Point-CMAE constructs contrastive input pairs by masking a given point cloud token twice randomly instead of applying he	✓	0.16
Point-CMAE uses a weight-sharing encoder and two identically structured decoders.	✓	0.18
Point-CMAE uses Chamfer distance loss as the reconstruction constraint.	×	0.06

References

- <http://arxiv.org/abs/2508.08910v1>
- <http://arxiv.org/abs/2407.05862v1>
- <http://arxiv.org/abs/2502.08347v1>