

Cross View CenterPoint for 3D Object Detection

Eshan Iyer
eshan@eshaniyer.tech

Dev Arora
aroradev314@gmail.com

Feng Chen
feng.chen@utdallas.edu

1 Introduction

The detection and localization of objects in 3D space play a crucial role in various applications, ranging from autonomous driving to robotics and augmented reality. Traditional methods for 3D object detection often rely on LiDAR sensors to capture point clouds, which provide accurate distance information but can be limited in capturing comprehensive scene details from multiple viewpoints. On the other hand, multi-camera systems can offer rich visual information from diverse perspectives but may need more precise depth information.

In this research, we propose a novel 3D object detection framework called Cross-View CenterPoint (CVCP) that leverages the strengths of multi-camera systems to enhance the accuracy and robustness of object detection. Our approach combines two powerful components: the Cross-View Transformers (CVT) as a backbone and the CenterPoint object detection model. By doing so, we aim to overcome the limitations of individual sensors and create a unified and more effective system for 3D object detection.

1.1 Background

LiDAR-based approaches widely adopt 3D object detection due to their precise depth measurements and robustness in various environmental conditions. However, point clouds' sparsity and inability to capture fine-grained visual details can limit these methods. In contrast, multi-camera systems can provide high-resolution images from different viewpoints, enabling a more comprehensive understanding of the scene. Using cameras could lead to a more complete and accurate representation of the environment. A traditional array of cameras can be significantly cheaper than a LiDAR system, allowing for greater accessibility.

1.2 Objectives

The primary objective of this research is to develop a novel 3D object detection framework that integrates Cross-View Transformers (CVT) and CenterPoint to leverage multi-view features for accurate and robust object localization. By using the information from camera images, we aim to achieve the following goals:

1. **Enhance Detection Accuracy:** By utilizing cross-view attention mechanisms in CVT, the proposed model can

capture and reason about correspondences across multiple camera views, leading to improved accuracy in object detection.

2. **Comprehensive Scene Understanding:** The incorporation of multi-camera features allows the model to benefit from rich visual cues, providing a more holistic understanding of the environment.
3. **Localization Precision:** The two-stage integration of CenterPoint refines object localization, resulting in more accurate and detailed 3D bounding box predictions.

1.3 Contributions

The contributions of this research are as follows:

1. **Cross-View CenterPoint (CVCP) Model:** We propose a novel architecture that combines Cross-View Transformers and CenterPoint to effectively integrate multi-view features into the 3D object detection process.
2. **Improved Object Detection Performance:** The CVCP model demonstrates superior performance compared to traditional LiDAR-based approaches by leveraging the complementary strengths of a multi-camera systems.
3. **Comprehensive Evaluation:** We conduct extensive experiments on benchmark datasets to validate the effectiveness and robustness of the proposed CVCP model in 3D object detection tasks.
4. **New Directions for Multi-View-Based Object Detection:** The integration of CVT and CenterPoint opens up new possibilities for multi-view-based object detection, contributing to advancements in computer vision and autonomous systems.

In the following sections, we provide a detailed overview of the Cross-View CenterPoint (CVCP) model, the integration of Cross-View Transformers, and the two-stage CenterPoint refinement. We present experimental results and analyses, comparing the performance of the proposed model with baseline approaches. Ultimately, we believe that the CVCP framework has the potential to significantly advance 3D object detection in real-world applications, where the fusion of multi-view features can lead to more accurate and comprehensive scene understanding.

2 Proposed Model: Cross-View CenterPoint (CVCP)

PLACEHOLDER

2.1 Architecture Overview

PLACEHOLDER

2.2 Cross-View Transformers as Backbone

PLACEHOLDER

2.2.1 Multi-View Input Representation

PLACEHOLDER

2.2.2 Cross-View Attention Mechanism

PLACEHOLDER

2.2.3 Camera-Aware Positional Encoding

PLACEHOLDER

2.2.4 Map-View Latent Embedding

PLACEHOLDER

2.3 CenterPoint Object Detection

PLACEHOLDER

2.3.1 Center Heatmap Head

PLACEHOLDER

2.3.2 Regression Heads

PLACEHOLDER

2.3.3 Velocity Head and Object Tracking

PLACEHOLDER

3 Two-Stage CenterPoint Integration

PLACEHOLDER

3.1 Second Stage with Point-Feature Extractor

PLACEHOLDER

3.1.1 Extracting Point-Features from Predicted Bounding Boxes

PLACEHOLDER

3.1.2 MLP for Confidence Score and Box Refinement

PLACEHOLDER

3.2 Class-Agnostic Confidence Score Prediction

PLACEHOLDER

3.3 Box Regression

PLACEHOLDER

4 Experimental Setup

PLACEHOLDER

4.1 Dataset Description

PLACEHOLDER

4.2 Implementation Details

PLACEHOLDER

4.3 Evaluation Metrics

PLACEHOLDER

5 Results and Analysis

PLACEHOLDER

5.1 Quantitative Results

PLACEHOLDER

5.2 Qualitative Results

PLACEHOLDER

5.3 Comparison with Baseline Models

PLACEHOLDER

6 Discussion

PLACEHOLDER

6.1 Advantages of CVCP

PLACEHOLDER

6.2 Limitations and Future Directions

PLACEHOLDER

7 Conclusion

PLACEHOLDER