# Transformers for 3D Object Detection in LiDAR Point Clouds for Autonomous Vehicles

GitHub repo: https://github.com/abubakar1107/Transformer_for_3d_obj_detection_in_LidarPC.git
NOTE: **Main** code is in the *model_test_v1* branch

1. Abubakar Siddiq                                                    [UID: 120403422]
2. Jangama Reddy Dhana Santhosh Reddy    [UID: 120405570]
3. Venkata Sai Charan Paladugu                          [UID: 120392948]

## Abstract

*The rise of autonomous vehicles demands highly accurate 3D object detection systems to ensure safe navigation in dynamic environments. LiDAR point clouds, with their sparse and irregular data structure, present significant challenges for existing detection frameworks. This project proposes a unique transformer-based architecture tailored for 3D object detection within LiDAR point clouds, leveraging self-attention mechanisms to enhance spatial correlation and long-range dependency capture. It uses pretrained PointNet model to produce feature embeddings and these embeddings are passed to the transformer blocks. By utilizing datasets like KITTI and employing appropriate data preprocessing, transformations using calibration files and positional encoding, the model aims to achieve high detection accuracy for vehicles, pedestrians, and cyclists in urban environments. The proposed solution is evaluated through metrics such as Intersection over Union (IoU) and cross-entropy loss, targeting improvements over state-of-the-art convolutional and PointNet based approaches.*

## 1. Introduction

In the era of autonomous vehicles, reliable perception systems are indispensable for ensuring safe navigation and decision-making in dynamic environments. LiDAR sensors, widely adopted in self-driving technologies, provide rich 3D point cloud data that capture spatial and structural information about the surroundings. This data is crucial for detecting and classifying objects such as vehicles, pedestrians, and cyclists. However, the unique characteristics of LiDAR are point cloud sparsity, irregularity, and non-uniform density pose significant challenges for existing detection algorithms.

To address these limitations, this project proposes a transformer-based 3D object detection framework for LiDAR point clouds. By leveraging self-attention mechanisms and pretrained PoinNet, the model aims to enhance detection accuracy and efficiency. Using datasets

such as KITTI, the model was evaluated on metrics like Intersection over Union (IoU) and classification accuracy. Unlike conventional CNN-based or PointNet architectures, our approach leverages the self-attention mechanism of transformers to capture global spatial relationships within the point cloud. This enables more effective modeling of long-range dependencies, crucial for accurate detection of objects at varying distances. Furthermore, we integrate a pre-trained PointNet++ module for robust feature extraction, eliminating the need to train a separate embedding network. Our results demonstrate that this combined approach enhances both detection accuracy and computational efficiency.

### 1.1. Problem Statement

Accurate 3D object detection in LiDAR point clouds is a critical requirement for autonomous vehicles to navigate safely and efficiently in complex, dynamic environments. Current detection frameworks, including Convolutional Neural Networks (CNNs) and PointNet-based architectures, face significant limitations:

1. Inadequate Long-Range Dependency Modeling: Existing methods struggle to capture global spatial relationships and dependencies in sparse and irregular LiDAR data, leading to suboptimal detection of objects at varying distances.

2. Computational Inefficiency: Many state-of-the-art models are computationally intensive, making them unsuitable for real-time applications required in autonomous driving.

3. Challenges with Data Sparsity: LiDAR point clouds exhibit uneven density, particularly in distant regions, which hampers the ability of conventional models to detect and localize objects accurately.

This project seeks to address these challenges by developing a transformer-based 3D object detection framework tailored for LiDAR point clouds. By leveraging self-attention mechanisms and integrating multi-modal data

sources, the proposed model aims to enhance detection accuracy, computational efficiency, and robustness in real-world urban driving scenarios.

## 1.2. Proposed Approach

This project addresses these challenges by implementing a transformer-based architecture, designed to leverage the self-attention mechanism's strength in capturing long-range dependencies and global spatial correlations. The approach uses a pretrained PointNet to produce features of the point cloud and these features are used as embeddings and passed into the transformer blocks. Key aspects of the project include:

- Preprocessing techniques for efficient handling of large-scale point cloud data. (Processing the point cloud in chunks).
- Custom loss functions combining smooth L1, Cross Entropy loss and IoU loss for precise 3D bounding box regression.
- Modular network designs evolve from basic embeddings to scalable transformer-based solutions.

## 2. Literature Review

Transformers, originally designed for NLP, are now employed for 3D object detection due to their ability to capture long-range dependencies in sparse data, such as LiDAR point clouds. The survey by Zhu et al. (2024) highlights the limitations of CNNs in 3D detection and demonstrates that transformer architectures like DETR can integrate local and global features effectively for complex scenarios but limited to image data.

Li3DeTr, a transformer-based approach, employs sparse convolutions and multi-scale attention to enhance 3D bounding box predictions. This method eliminates the need for non-maximum suppression, achieving state-of-the-art results on datasets like nuScenes. Similarly, PolarStream uses a polar grid representation for computational efficiency, improving latency without sacrificing accuracy.

Corral-Soto et al. (2023) demonstrate that optimizing the density of point clouds across distance ranges can enhance detection without architectural changes. These complements methods using data augmentation to enrich training datasets.

Models like FPA-3DOD address the latency issue inherent in 3D detection by utilizing polar space representations and lightweight architectures, crucial for real-time applications in autonomous driving.

### 2.1. Challenges with Sparse Data

Point cloud sparsity and irregularity pose challenges for traditional 3D detection approaches. Graph Neural Networks (GNNs) have been explored as an alternative to address these issues by operating on graph-structured data.

### 2.2. Research Gaps and Areas for Improvement

1. Multi-Modal Integration: Existing models struggle with seamless integration of LiDAR and RGB data in transformer architectures.

2. Scalability: Ensuring real-time performance on diverse datasets like KITTI and Waymo remains challenging.

3. Data Sparsity: Approaches to handle sparsity in long-range detections, especially in urban scenarios, need refinement.

While Li3DeTr [11] effectively utilizes sparse convolutions and multi-scale attention, its focus is primarily on improving bounding box predictions. In contrast, our approach emphasizes the integration of a pre-trained PointNet++ backbone to enhance feature learning. Furthermore, while PolarStream [13] addresses computational efficiency using a polar grid representation, its performance on complex urban scenes with varying point densities remains to be explored. Our work directly tackles this challenge by incorporating positional encoding and a chunking mechanism to handle large-scale, unevenly distributed point clouds.

## 3. Dataset

The KITTI dataset was selected for this project, as the Waymo Open Dataset posed challenges related to size and storage accessibility in the public GCS bucket. KITTI dataset includes training and testing folders each contain multimodal data of *Velodyne point cloud* and *RGB images* with labels containing object bounding boxes in 3D and 2D along with a classification label. Raw Velodyne point cloud data from the KITTI dataset underwent preprocessing, including:

1. Range-based filtering: Points beyond a 80-meter radius were removed to focus on relevant objects.

2. Normalization: Point coordinates were normalized to a range of [0, 1] to ensure consistent scaling across different scenes.

## 3.1. Input/output

Input: LiDAR point cloud data representing a 3D spatial view of the environment.

Output: The model will generate 3D bounding box predictions for objects in the scene, along with classification labels.

In the context of training an object detection model using the KITTI dataset, it is crucial to transform raw Velodyne point cloud data into a format suitable for neural network input. This process involves several steps, including loading the raw data, preprocessing it, and converting it into tensors that can be fed into the model.

The KITTI dataset is a widely used benchmark dataset for autonomous driving research. It contains various types of data collected from a moving vehicle, including:

1. Images: High-resolution color and grayscale images captured by cameras mounted on the vehicle.

2. Annotations: Ground truth labels for objects in the images, including their class and bounding box coordinates.

3. Calibration Files: Parameters for mapping 3D world coordinates to the 2D image plane.

4. Velodyne Point Clouds: 3D point cloud data captured by a Velodyne LiDAR sensor.

The initial step involves loading the raw Velodyne point cloud data. The data is typically stored in binary files, where each point is represented by its coordinates (x, y, z) and intensity. This raw data is read and reshaped into a structured format, such as a 2D array, where each row corresponds to a single point.

Once the raw data is loaded, it undergoes a series of preprocessing steps. These steps are essential to ensure the data is in a consistent and usable state for the model. Preprocessing may include filtering out points based on specific criteria, such as removing points with certain coordinate values that are deemed irrelevant or noisy. Additionally, normalization techniques may be applied to

standardize the coordinates and intensity values, ensuring that the data falls within a specific range.

After preprocessing, the data is converted into tensors. Tensors are the primary data structure used in PyTorch for model input. This conversion is necessary because neural networks in PyTorch are designed to operate on tensor data. The preprocessed point cloud data is transformed into a tensor format, which involves specifying the data type and ensuring the tensor dimensions align with the model's expected input shape.

The entire transformation process from raw data to model input involves orchestrating the loading, preprocessing, and tensor conversion steps. This systematic approach ensures that the data is properly prepared and standardized, making it suitable for training the object detection model. The raw Velodyne point cloud data is effectively transformed into a format that the neural network can process, enabling the model to learn from the input data and make accurate predictions.

## 3.2. Evaluation Metrics

Model performance is evaluated using Intersection over Union (IoU) for bounding box overlap, cross-entropy loss for classification. This will ultimately evaluate the performance of the model on how well it is able to capture the point cloud space of the detected objects and label it accordingly.

The loss function combines multiple components to train the model effectively. The regression loss utilizes Smooth L1 to measure the difference between predicted and ground truth bounding box parameters, ensuring robustness against outliers. This is augmented with an IoU-based term that accounts for the overlap between predicted and ground truth boxes, providing a direct measure of localization quality. Together, these terms capture both the scale and positioning accuracy of the detected objects.

Additionally, the classification loss applies a weighted cross-entropy function to handle class imbalances in the dataset. By integrating these three terms, the loss function balances the need for precise object localization with accurate class prediction. The weights of the components are tunable via hyperparameters, allowing flexibility to prioritize certain objectives depending on the task requirements.

Classification Loss:

We utilize weighted cross-entropy loss for object classification. This loss function effectively penalizes incorrect class predictions, and the weighting factor allows us to address potential class imbalances in the training data. Formally, the classification loss is defined as:

$$L_{\text{classification}} = -\sum_c w_c \cdot y_c \cdot \log(\hat{y}_c)$$

Regression Loss:

The bounding box regression, we employ the Smooth L1 loss function. This loss function is less sensitive to outliers compared to the L2 loss (mean squared error), providing more robust training. The Smooth L1 loss is defined as:

$$L_{\text{regression}} = \begin{cases} 0.5 \cdot (p - t)^2 & \text{if } |p - t| < 1 \\ |p - t| - 0.5 & \text{otherwise} \end{cases}$$

Intersection over Union (IoU)

This is a crucial metric in object detection, measuring how well a predicted bounding box overlaps with the actual (ground truth) bounding box of an object.

Intersection volume:

This calculates the volume of the overlapping region between the predicted and ground truth boxes.

$$\text{Intersection} = \prod_{d=1}^{3} max\big(0, min(p_{d+3}, t_{d+3}) - max(p_d, t_d)\big)$$

Where $p_d$ and $t_d$ are the coordinates of the bounding box.

Predicted and true volumes:

These calculate the volumes of the individual predicted and ground truth bounding boxes.

$$\text{Volume}_{\text{pred}} = \prod_{d=1}^{3} (p_{d+3} - p_d)$$

$$\text{Volume}_{\text{true}} = \prod_{d=1}^{3} (t_{d+3} - t_d)$$

IoU (Intersection over Union):

This is the ratio of the intersection volume to the union volume. The union volume is the total volume encompassed by both boxes combined, minus the overlap (to avoid double-counting).

$$\text{IoU} = \frac{\text{Intersection}}{\text{Volume}_{\text{pred}} + \text{Volume}_{\text{true}} - \text{Intersection} + \epsilon}$$

IoU (Intersection over Union)Loss:

This measures how far the IoU is from the ideal value of 1 (perfect overlap).

$$IoU\ Loss \quad L_{\text{IoU}} = 1 - \text{IoU}$$

Combined Loss:

This combines the IoU loss with other losses, likely a regression loss for predicting the exact bounding box coordinates and a classification loss for predicting the object's class. The total combined loss is given by:

$$L_{\text{total}} = \alpha \cdot L_{\text{regression}} + \beta \cdot L_{\text{IoU}} + L_{\text{classification}}$$

Where:

- $\alpha$: Weight for regression loss
- $\beta$: Weight for IoU loss

## 4. Model Architecture

Unlike convolutional-based models, our approach uniquely uses transformer-based architectures with multi-headed attention mechanisms to capture long-range dependencies in the point cloud more effectively and the embeddings to this transformer blocks is obtained by a pretrained PointNet++ model. This framework aims to process large-scale point cloud data more efficiently and accurately and ultimately produces a better output in detection and classification of the objects in the 3D point cloud. Our model is peculiar is it uses of a pretrained PointNet++ as the embedding block to produce initial embeddings which will be further utilized to perform attention on them. By using this we can get meaningful embeddings to be passed to the transformer blocks so that

we don't need to train another network to produce the embeddings.

### 4.1. Technical Approach:

The technical approach includes several key steps:

**Data Processing**: LiDAR data is filtered and calibrated using provided calibration files in the dataset. Data is processed in chunks, with each chunk containing 5,000 points to improve training efficiency. Once a chunk is processed the output features are stored and the next chunk is processed to add to the previous features till a particular lidar file is completed.

**Model Details**: Various networks have been designed for testing:

*Network1*: A transformer-based model using multi-headed attention without positional encoding, featuring a straightforward linear embedding. The embeddings are very trivial (N, 3) to (N, 64) where 64 is the final embedding length of each point in the point cloud. This linear embedding is not enough because this can't capture the higher dimensional features.

The smaller model summary is as follows:

```
==================================================================
Layer (type:depth-idx)                   Output Shape         Param #
==================================================================
├─get_model: 1-1                         [-1, 512, 64]        --
│    └─PointNetSetAbstractionMsg: 2-1    [-1, 3, 512]         --
│    └─ConvId: 2-2                       [-1, 64, 512]        (20,544)
├─Linear: 1-2                            [-1, 512, 64]        4,160
├─ModuleList: 1                          []                   --
│    └─PointTransformerBlockWithPE: 2-3  [-1, 512, 64]        --
│    │    └─PositionalEncoding: 3-1      [-1, 512, 64]        --
│    │    └─MultiheadAttention: 3-2      [-1, 1, 64]          16,640
│    │    └─LayerNorm: 3-3               [-1, 512, 64]        128
│    │    └─Sequential: 3-4              [-1, 512, 64]        8,320
│    │    └─LayerNorm: 3-5               [-1, 512, 64]        128
│    └─PointTransformerBlockWithPE: 2-4  [-1, 512, 64]        --
│    │    └─PositionalEncoding: 3-6      [-1, 512, 64]        --
│    │    └─MultiheadAttention: 3-7      [-1, 1, 64]          16,640
│    │    └─LayerNorm: 3-8               [-1, 512, 64]        128
│    │    └─Sequential: 3-9              [-1, 512, 64]        8,320
│    │    └─LayerNorm: 3-10              [-1, 512, 64]        128
│    └─PointTransformerBlockWithPE: 2-5  [-1, 512, 64]        --
│    │    └─PositionalEncoding: 3-11     [-1, 512, 64]        --
│    │    └─MultiheadAttention: 3-12     [-1, 1, 64]          16,640
│    │    └─LayerNorm: 3-13              [-1, 512, 64]        128
│    │    └─Sequential: 3-14             [-1, 512, 64]        8,320
│    │    └─LayerNorm: 3-15              [-1, 512, 64]        128
│    └─PointTransformerBlockWithPE: 2-6  [-1, 512, 64]        --
│    │    └─PositionalEncoding: 3-16     [-1, 512, 64]        --
│    │    └─MultiheadAttention: 3-17     [-1, 1, 64]          16,640
│    │    └─LayerNorm: 3-18              [-1, 512, 64]        128
│    │    └─Sequential: 3-19             [-1, 512, 64]        8,320
│    │    └─LayerNorm: 3-20              [-1, 512, 64]        128
├─Sequential: 1-3                        [-1, 3]              --
│    └─Linear: 2-7                       [-1, 128]            8,320
│    └─ReLU: 2-8                         [-1, 128]            --
│    └─Linear: 2-9                       [-1, 3]              387
├─Sequential: 1-4                        [-1, 7]              --
│    └─Linear: 2-10                      [-1, 128]            8,320
│    └─ReLU: 2-11                        [-1, 128]            --
│    └─Linear: 2-12                      [-1, 7]              903
├─Sequential: 1-5                        [-1, 7]              (recursive)
│    └─Linear: 2-13                      [-1, 128]            (recursive)
│    └─ReLU: 2-14                        [-1, 128]            --
│    └─Linear: 2-15                      [-1, 7]              (recursive)
==================================================================
Total params: 143,898
Trainable params: 122,954
Non-trainable params: 20,544
Total mult-adds (M): 11.08
```

*Network2*: This builds on *Network1* by incorporating positional encoding to capture the spatial relations between each point in a lidar file.

*Network3*: This builds on *Network2* and uses a pretrained PointNet++ model as an embedding block to extract complex features. PointNet++ processes the input features (N, 3) and produces a (N, 64) feature space which is relevant to the Lidar file, PointNet++ is designed to capture localized features and normal in a point cloud and produces a meaningful embedding. These embeddings are then processed through transformer blocks, yielding high-dimensional features which will be passed to the fully connected layers and transformed to a much higher space and finally these outputs are passed to a classification head and a regression head which will produce output bounding box with a classification label.
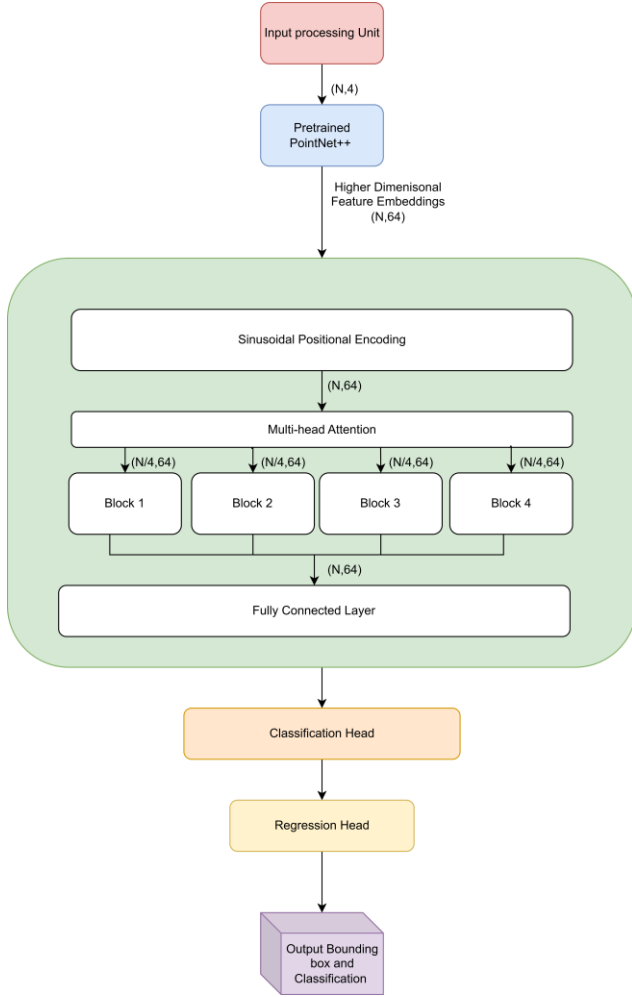
*Network4*: Till now, the network only has the capability to process one fixed chunk of points but, generally this is not enough, as each point cloud contains more than 5000 points, so batch wise processing is necessary, and it is incorporated into this network. This processes data in chunks, enabling the model to handle the large sequence length of point cloud files.

*Network5*: Integrated insights from all earlier networks and made it scalable so that it can be made complex to capture better features from the data and perform better on the validation dataset. This network utilizes a custom loss function made by combining L1 loss, IoU loss, and cross-entropy loss to improve object localization and classification.

### 4.2 Network Architecture

The *Network 5* is finalized as the final model begins by leveraging a pre-trained PointNet, which is responsible for extracting meaningful geometric features from the raw input points. This backbone effectively captures local and global patterns in the data. To ensure these features align with subsequent processing requirements, a feature dimension adaptation layer refines the extracted embeddings, transforming the output from 512 dimensions to a 64-dimensional space.

(The flowchart represents the model architecture we designed, the number of blocks and heads was experimentally determined later, the number of heads we see in the below picture is an example model configuration with 4 blocks of multiheaded attention each with 4 heads. The flow chart is of one Transformer block with four heads)

A series of four transformer blocks follow, each equipped with four attention heads and positional encoding to integrate spatial context into the embeddings. These blocks (heads in this picture) utilize self-attention mechanisms, allowing the model to dynamically focus on different regions of the point cloud. After the attention layers, the data retains its dimensionality of 64 but incorporates enhanced spatial relationships. The inclusion of feed-forward layers and normalization steps enhances the model's ability to learn complex patterns while maintaining stability during training.

To manage large-scale point clouds, the architecture includes a chunking mechanism that processes data in manageable segments. This approach ensures scalability without compromising performance. Features extracted from individual chunks are aggregated using mean pooling, providing a holistic representation of the input point cloud.

Finally, the model follows two distinct heads. The classification head predicts the object class by processing the aggregated features through dense layers, reducing the 64-dimensional embedding to 3 output classes. Meanwhile, the regression head estimates precise bounding box parameters, including position, dimensions, and orientation, producing 7 outputs corresponding to x, y, z, width, height, depth, and yaw. Both heads work in tandem to deliver robust object detection outputs, capable of distinguishing objects and localizing them accurately in 3D space.

Overall, our architecture is a well-balanced design, combining pre-trained feature extraction, transformer-based spatial reasoning, and tailored mechanisms for large-scale data handling. Its modularity and efficiency make it suitable for real-world 3D object detection tasks.

## 5. Training and Results

Several experiments were conducted to explore the impact of different hyperparameters on model performance, the GPU requirement for such an architecture and dataset size is tremendously huge and we did the training on remote cloud instances that we bought in *Lambda Labs*. We experimented with different configurations and finalized the best model.

### 5.1. Challenges in Achieving Lower Loss

In the process of training the transformer model, the primary objective is to minimize the loss effectively. However, despite rigorous hyperparameter tuning and optimization, the loss values remained higher than expected. This suggests that additional factors, such as computational resource constraints, might play a significant role in improving the model's performance.

### 5.2. Hyperparameter Tuning Experiments

- *Number of heads in multi-headed attention*: Multi-headed attention was tested with varying numbers of heads to balance the model's capacity to learn diverse relationships while maintaining computational feasibility. A higher number of attention heads allows the model to focus on multiple parts of the input sequence simultaneously, enhancing the model's ability to capture intricate dependencies. However, an increase in the number of heads also adds computational overhead and memory usage.
- *Number of Points processed in Chunks:* Each head processes a subset of the input points independently. With multi-headed attention, the number of points processed at the same time is influenced by the embedding dimension and the number of heads.
- *Learning-Rate:* Fine-tuning the learning rate was critical in balancing fast convergence with stability.

Lower learning rates (e.g., 1e-4) offered more stable training trajectories, while slightly higher rates (e.g., 5e-4) showed faster convergence but occasionally led to instability. Experimentation with learning rate schedulers, such as cosine annealing and warm-up phases, provided additional improvements in optimization.

- *Embedding Dimension:* Embedding dimensions of 64 and 128 were tested to balance computational efficiency and feature representation. 64 for resource-constrained environments, offering efficient training but with potentially limited expressiveness. 128 Provided richer feature representations, enabling the model to learn more intricate data patterns. However, this required significantly more memory and computation.

- *Number of transformer Blocks:* Increasing the number of attention blocks was tested, as deeper models often capture more intricate patterns in data. However, the loss improved only marginally due to possible underfitting caused by other factors.

- *Batch Size Variation*: Both smaller and larger batch sizes were tested. While smaller batches allowed the model to train on more diverse subsets of data, larger batches provided a more stable optimization trajectory. Neither approach led to a significant decrease in the loss.

### 5.3. Observations on Resource Constraints

The experiments revealed that the GPU resources available were a limiting factor. High-end GPUs typically enable:

- Larger batch sizes, reducing gradient noise and improving optimization.
- Faster convergence due to increased computational power.
- Support for deeper architectures without memory bottlenecks.

The limited resources restricted the ability to fully explore these enhancements.

### 4. Loss Trends

The loss curves below illustrate the training loss over iterations, highlighting the impact of hyperparameter changes. While the model shows a downward trend in loss, it plateaus early, suggesting room for improvement with better resources.
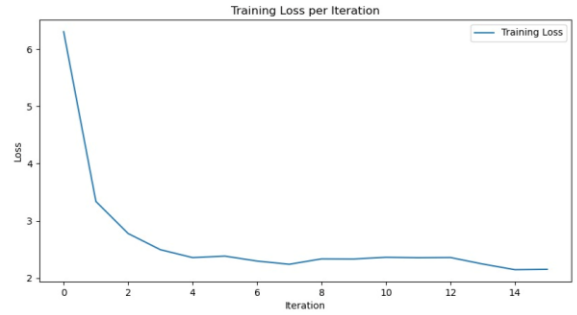


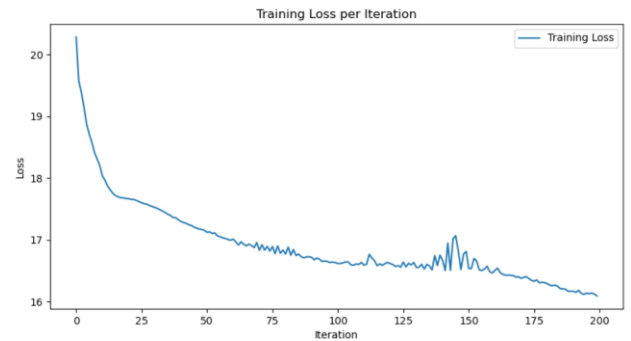*Figure 1: Training Loss over Iteration for small dataset size and larger batches for complex model.*



*Figure 2: Training Loss over Iteration for Large dataset size and larger batches for complex model.*
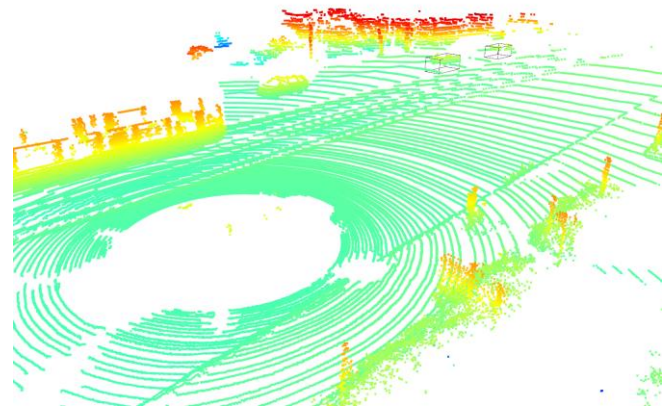


*Figure 3: Example Detection using best_model_saved.pth*

### 6. Future Work

To further improve performance:

- Access to high-performance GPUs is recommended for training with larger batch sizes and deeper architectures.
- Advanced hyperparameter optimization methods,

such as Bayesian optimization or grid search, could be explored.

- Fine-tuning the learning rate and experimenting with alternative regularization methods may also help.
- Adjusting the batch size, number of blocks, attention heads with more computational power could give better results with lower loss and greater accuracy.

Through extensive experiments on the KITTI dataset, we show that our transformer-based model achieves competitive performance in 3D object detection, particularly for challenging scenarios with sparse point clouds. Our results highlight the effectiveness of positional encoding in capturing spatial context and the importance of efficient data chunking for handling large-scale point cloud data.

## 7. Conclusion

The transformer-based architecture for 3D object detection in LiDAR point clouds demonstrates promising results while highlighting several key achievements and challenges. The integration of pre-trained PointNet++ with transformer blocks effectively captures both local geometric features and global spatial relationships in point cloud data. The model's chunking mechanism successfully addresses the challenge of processing large-scale point clouds, enabling scalable performance without compromising accuracy. Despite resource constraints limiting full exploration of deeper architectures, the model shows competitive performance on the KITTI dataset, particularly in handling sparse point clouds and capturing spatial context through positional encoding. The custom loss function combining L1, IoU, and cross-entropy losses proves effective for joint optimization of object localization and classification tasks.

**References**

[1] KITTI Dataset: Geiger, A., Lenz, P., & Urtasun, R. (n.d.). Vision meets robotics: The KITTI dataset. Retrieved from https://www.cvlibs.net/datasets/kitti/

[2] Zhao, H., Jiang, L., Jia, J., Torralba, A., & Koltun, V. (2020). Point Transformer. arXiv. Retrieved from https://arxiv.org/abs/2012.09164.

[3] Arzhanov, A. (2019). *3D Object Detection from Point Cloud*. CS230 Deep Learning, Stanford University. Retrieved from your provided file.

[4] Zhu, M., Gong, Y., Tian, C., & Zhu, Z. (2024). A systematic survey of transformer-based 3D object detection for autonomous driving: Methods, challenges, and trends. *Drones, 8*(412). https://doi.org/10.3390/drones8080412

[5] Raut, G., & Patole, A. (n.d.). *End-to-End 3D Object Detection Using LiDAR Point Cloud*. University of Maryland. Retrieved from your provided file.

[6] Bhat, M., Han, S., & Porikli, F. (2021). Fast polar attentive 3D object detection on LiDAR point clouds. ("Fast Polar Attentive 3D Object Detection on LiDAR Point Clouds") *Neural Information Processing Systems (NeurIPS) Workshop*. Retrieved from your provided file.

[7] Corral-Soto, E. R., Grandhi, A., He, Y. Y., Rochan, M., & Liu, B. (2023). "Improving LiDAR 3D object detection via range-based point cloud density optimization." ("Improving LiDAR 3D Object Detection via Range-based Point ... - Synthical") *Huawei Noah's Ark Lab*. Retrieved from your provided file.

[8] Li, X., Wang, C., Wang, S., Zeng, Z., Liu, J., & Zheng, B. (2024). Improving the point cloud-based 3D object detection for autonomous driving by constructing multi-scale features. *Chongqing University and Zhejiang University*. Retrieved from your provided file.

[9] Erabati, G. K., & Araujo, H. (2022). *Li3DeTr: A LiDAR-Based 3D Detection Transformer*. Institute of Systems and Robotics, University of Coimbra. Retrieved from your provided file.

[10] Chen, Q., Vora, S., & Beijbom, O. (2021). "PolarStream: Streaming LiDAR object detection and segmentation with polar pillars." ("PolarStream | Proceedings of the 35th International Conference on ...") *Neural Information Processing Systems (NeurIPS) Workshop*. Retrieved from your provided file.