

CloudAttention: Efficient Multi-Scale Attention Scheme For 3D Point Cloud Learning

Mahdi Saleh^{1*} Yige Wang^{1*} Nassir Navab² Benjamin Busam¹ Federico Tombari³

Abstract—Processing 3D data efficiently has always been a challenge. Spatial operations on large-scale point clouds, stored as sparse data, require extra cost. Attracted by the success of transformers, researchers are using multi-head attention for vision tasks. However, attention calculations in transformers come with quadratic complexity in the number of inputs and miss spatial intuition on sets like point clouds. We redesign set transformers in this work and incorporate them into a hierarchical framework for shape classification and part and scene segmentation. We propose our local attention unit, which captures features in a spatial neighborhood. We also compute efficient and dynamic global cross attentions by leveraging sampling and grouping at each iteration. Finally, to mitigate the non-heterogeneity of point clouds, we propose an efficient Multi-Scale Tokenization (MST), which extracts scale-invariant tokens for attention operations. The proposed hierarchical model achieves state-of-the-art shape classification in mean accuracy and yields results on par with the previous segmentation methods while requiring significantly fewer computations. Our proposed architecture predicts segmentation labels with around half the latency and parameter count of the previous most efficient method with comparable performance. The code is available at <https://github.com/YigeWang-WHU/CloudAttention>.

I. INTRODUCTION

Point clouds are prevalent representations of 3D environments and objects which come naturally from range sensors. Point clouds are input to many robotics and augmented reality applications. Despite their wide usability, point clouds have multiple drawbacks, such as sparsity and density variations. 3D scans usually consist of millions of points represented sparsely in computers. Processing point clouds invariant to their permutation is crucial in any point cloud learning pipeline. At the same time, applying order-invariant and symmetric functions on large-scale point clouds is expensive when subtle underlying structures are relevant. Recent works suggest iterative sampling and applying order-invariant operations to point groups.

For instance, PointNet++ [1] suggests grouping points using a ball query and applies vanilla PointNet [2] on

*the authors contributed equally to this paper

**This work was partly sponsored by the German Federal Ministry for Economic Affairs and Energy (grant number: 19A19005B) through the VDA KI-Absicherung project.

¹Mahdi Saleh, Yige Wang and Benjamin Busam are with the Faculty of Computer Science, Technische Universität München (TUM), 85748 Garching bei München, Germany m.saleh@tum.de, yige.wang@tum.de, b.busam@tum.de

²Nassir Navab is with the Faculty of Computer Science, Technische Universität München (TUM), 85748 Garching bei München, Germany, and also with the Department of Computer Science, Johns Hopkins University, Baltimore, 21218 MD USA nassir.navab@tum.de

³Federico Tombari is with the Faculty of Computer Science, Technische Universität München (TUM), 85748 Garching bei München, Germany, and also with Google, 8002 Zurich, Switzerland tombari@in.tum.de

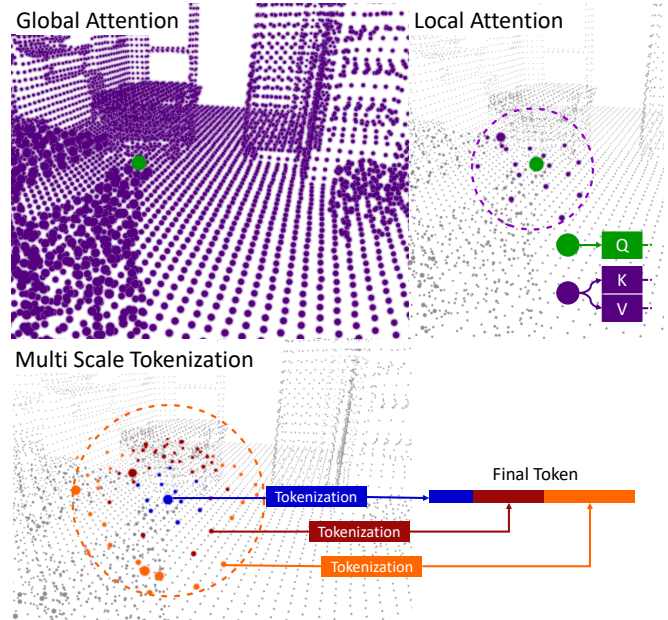


Fig. 1. Top: We extract global and local attention using multi-head attention concepts, with Query Q, Key K, and Value V. Bottom: We propose an efficient multi-scale tokenization strategy to describe non-heterogeneous point clouds effectively.

each group followed by set abstraction. While there is a hierarchical notion to PointNet++, it only focuses on the fix-sized receptive field provided by the ball query at each level. Following methods also utilize a fixed grouping using K-NN or ball query [3], [4]. Another issue with PointNet-based methods is the extensive use of max-pooling operations, through which the gradient flow is broken and local relationships are lost.

Transformers [5] have shown promising results in Natural Language Processing tasks and vision Transformers [6] show the potential of utilizing Multi-head attention units for image recognition [6]. Transformers are a good fit for sequential data since they calculate attention per node and use it for weighted pooling. Therefore, the gradient can flow to respective local neighbourhoods influenced by relationships and geometry. Transformers are applied to point sets [7], [8], [9] in previous works. PointTransformer [7] augments all point features by using a multi-head self-attention. In real-time use, where a huge point cloud exists, this comes with a high cost since this has quadratic complexity with respect to the number of points.

As an essential part of Transformers [6], self-attention is calculated to define self similarities. Using this, the nodes

that are featured similarly have greater attention and are weighted higher. In previous works in point cloud processing [7] self-attention units are used to construct self-similarities across the whole input data. In this work, we conduct self-attention only to a local receptive field to focus on self-similarities with keys in the neighboring points. This way, we not only save many computations, but also can guarantee a well-distributed local attention score.

Apart from self-attention, cross-attention is designed to capture inter-sequence attention. It is typically calculated from one query to all keys on the corresponding sequence. We use cross attention to capture global attention and feature similarities in our single frame setting. Inspired by set transformer [8], we use raw input point cloud as the cross-domain to find global similarities to every point. By tokenizing the point cloud data and reducing the number of query tokens iteratively, we significantly reduce the number of cross-attention calculations. We call our proposed iterative cross-attention module, Global Attention Unit (GAU). In several iterations, we apply a Local Attention Unit (LAU) alternately with our GAU unit. This way, we capture local to global attention through the feature extraction stage. Figure 1:top illustrates global and local attention.

Point cloud sparsity typically varies in each sample volume. Previous methods normally perform tokenization and abstraction in a single ball or use KNN [10], [11]. The ball query helps to keep the metric [1] but faces an in-balance in terms of a number of queried points. On the other hand, KNN forces the use of a fixed number of points in the neighborhood and loses metrics in varying densities [1]. This work also introduces an efficient hybrid approach where we use KNN with varying K within a ball to describe a region. Our multi-scale tokenization captures features equivariant to their level of sparsity and scale.

Combining the proposed modules reduces the computational complexity compared to previous methods and maintains a hierarchical understanding of the point cloud. The proposed efficient multi-scale feature encoder can be used in different downstream tasks such as classification and segmentation.

In summary, our contributions are the following:

- We propose an efficient hierarchical architecture to learn on point clouds using iterative attention.
- We define self- and cross-attention modules for large-scale point clouds incorporating spatial local, and global attention mechanisms.
- We propose an efficient Multi-Scale Tokenization (MST) which deals with the problem of heterogeneity of sparse point clouds.

II. RELATED WORKS

A. Point Cloud Networks

3D point clouds are unordered sparse data, and conventional CNN methods can not directly process them. Some methods quantize point clouds into regular grids and apply 3D CNNs on them [12], [13]. However, such voxel-based

approaches are computationally intensive and result in high latency.

Deep learning techniques can also process unordered sets directly using MLPs. PointNet [2] builds a permutation-invariant architecture based on pointwise MLP operators. Features are pooled to form a global high-dimensional feature vector representing the whole set. PointNet++ [1] extends PointNet by adding a multi-layer hierarchical procedure using furthest point sampling in each layer and grouping points by ball-query. These hierarchical operations are computationally expensive and can cause high runtime. Therefore, in large-scale applications, PointNet is usually combined with voxel-based methods [14], [13]

While previous methods project MLP features into high-dimensional space with loose geometric intuition, graph neural networks provide a framework to extract richer local features [4], [15]. PointCNN [10] uses spatial relationship for feature aggregation of neighboring points. EdgeConv [4] constructs graphs dynamically at each layer using K-NN in the embedding space. Earlier graph neural networks for point clouds are limited to node-level message passing and are impractical for large-scale point cloud applications. Following works [16], [5] combine node interactions and updates with edge-based attentions.

B. Attention and Transformers

Attention can be calculated in various forms to represent a weighting for every point or node. It can help focus on the most relevant parts of the input. Graph Attention Networks [16] use masked self-attentional layers to predict attention coefficients per edge, and a weighted sum is used to aggregate features per node in the graph.

In the transformer network [5] the authors propose the multi-head attention concept, which constitutes key, query, and value heads to calculate the output using scaled dot-product. The specialty of Transformers is that it is entirely based on attention operations, without convolutions or recurrent operations. Besides advances in machine translation, Transformer is applied to 2D [6], 3D[17], [18] and spatio-temporal [19] computer vision problems. Set transformer [8] proposes a self-attention mechanism for interactions among elements of a set based on a sparse Gaussian process. Following this, Perceiver [20] adds cross-attention to lower-dimensional data to reduce the complexity from quadratic to linear. Our GAU design is inspired by the cross-attentions in Perceiver.

Several works have applied transformers and self-attentions to 3D point clouds. PCT [21] is the first work to adapt the self-attention of Transformer to point cloud data. It incorporates Graph Laplacian and element-wise subtraction into the self-attention calculations. Point Transformer [11] applies the self-attentions for K-NNs and performs a global average pooling to get the global feature vector. In another Point Transformer work [7] local features are sorted by a proposed SortNet, and self-attentions are conducted to all other points for feature augmentations. While these works suggest

different self-attention concepts, they have shortcomings in costly global attention calculations.

We redesign self and cross attentions in an iterative architecture and use multi-scale tokenization to productively describe regions while saving on attention computations.

III. METHODOLOGY

This section explains our proposed method and modules one by one. Our CloudAttention extracts features from a sparse point cloud iteratively by applying local (LAU) and global attention (GAU) as well as multi-scale tokenization (MST) modules. Fig. 2 provides an overview of our pipeline for the semantic segmentation task.

A. Notation & Background

A point cloud is a set of points, and each point is represented by its 3D coordinate and optional features, for example, color features or surface normals. Formally, given an input point cloud with N points, we describe it as $P = \{(p_i, f_i)\}_{i=1}^N \in \mathbb{R}^{N \times d}$ where p_i is the 3D spatial coordinate of the i_{th} point and f_i its corresponding features. We stack coordinates and features vertically to form a matrix notation, i.e. $\mathbf{p} \in \mathbb{R}^{N \times 3}$ and $\mathbf{f} \in \mathbb{R}^{N \times (d-3)}$. For a classification task, the output is category C of the whole point cloud P , and for a segmentation task, a categorical label per every point in P .

The Multi-Head Attention (MHA) proposed in Transformer [5], is an explicit form of learnable attention. It projects the input into multiple query-key-value tuples and applies the attention within each set. The attentive values across sets are aggregated and applied with a linear layer.

B. Local Attention Unit (LAU)

Regular multi-head attention, is a global operation that models the relationship among all elements. However, it is not always necessary to relate elements globally. The majority of elements have negligible global relation to one another, especially regarding their low-level geometrical information. Moreover, learning global relationships in a huge point cloud is a redundant and expensive operation. Therefore, we propose a local self-attention unit (LAU) to model relationships within a vicinity. LAU is built upon self-attention concepts but only operates on a small number of neighbors. Formally, it is described as

$$\begin{aligned} \text{LAU} &= S + \text{FF}(S) \\ S &= X + \text{MHA}(X_i W_Q, X_j W_K, X_j W_V) \end{aligned} \quad (1)$$

where FF is simple feed-forward layer, $X_j \in \Omega(X_i, K)$ with $\Omega(X_i, K)$ is the set of K -nearest neighbors around X_i with respect to their Euclidean distance $\|X_i - X_j\|_2$.

As LAU only operates sparsely and on a small subset of tokens instead of all tokens, the computation cost is significantly reduced compared to previous dense self-attentions strategies.

C. Global Attention Unit (GAU)

Restricting attention and information exchange to small regions results in a loss of the global context, which is essential for scene segmentation tasks. Complementing our local attention, we also define a global attention function. We design a global attention unit (GAU) to incorporate global context and bring hierarchical understanding. The global attention unit is a cross-attention operation between tokens and the raw points.

$$\begin{aligned} \text{GAU} &= C + \text{FF}(C) \\ C &= T + \text{MHA}(T W_Q, P W_K, P W_V) \end{aligned} \quad (2)$$

where T denotes tokens and P represents raw input point clouds. Since GAU queries from a small number of tokens to all points, the computations do not scale quadratically with the entire number of input points. Furthermore, the number of tokens is reduced iteratively in the pipeline, enabling large-scale point cloud encoding with compact operations.

D. Multi-scale Tokenization (MST)

Non-uniformity is intrinsic in point clouds and imposes challenges to deep learning models. Point-based learning methods largely depend on sampling densities and sparsity of the input point clouds. Nevertheless, feature estimation and tokenization in a region should depend on the geometry rather than the point density. Thus, models that counteract the influences of varying densities at low cost are more applicable to real-world settings. To this end, we propose the multi-scale tokenization (MST) technique to allow the efficient processing of input point sets.

We initially sample several centroids using farthest point sampling (FPS)[1] to assure samples are distributed uniformly over the whole set. A ball query is used for each centroid to find relevant points in its vicinity. We then sort the queried points with respect to their Euclidean distance and construct a multi-scale feature within quantized distances. Figure 1 visualizes the MST operation on a point cloud region. Formally, we do multi-scale tokenization as follows:

$$\begin{aligned} \text{MST} &= T_{s_1} \oplus T_{s_2} \oplus \dots T_{s_N} \\ T_{s_i} &= \Theta_{p_j \in \Omega(c_i, K_i)}(\delta((p_j - c_i) \oplus f_j)) \end{aligned} \quad (3)$$

where N is the total number of scales, T_{s_i} is the feature abstracted at level s_i , the centroid coordinate is depicted as c_i and p_j belongs to its nearest K_i point set $\Omega(c_i, K_i)$. The operators δ describe a linear layer and Θ is a pooling operation, while \oplus denotes concatenation. We set $K_1 < K_2 < \dots < K_N$ to construct features at increasing scales.

MST is efficient for two reasons: 1) We use a single ball-query operation. We sort the points in the queried ball and simply form multi-scale features by tensor slicing, which omits any additional K -NN operations. 2) The linear layer δ is shared across scales, which avoids multiple MLP usage for features abstraction.

Using a hybrid multi-scale approach, we borrow the advantages of both ball-query and K -NN methods. The method

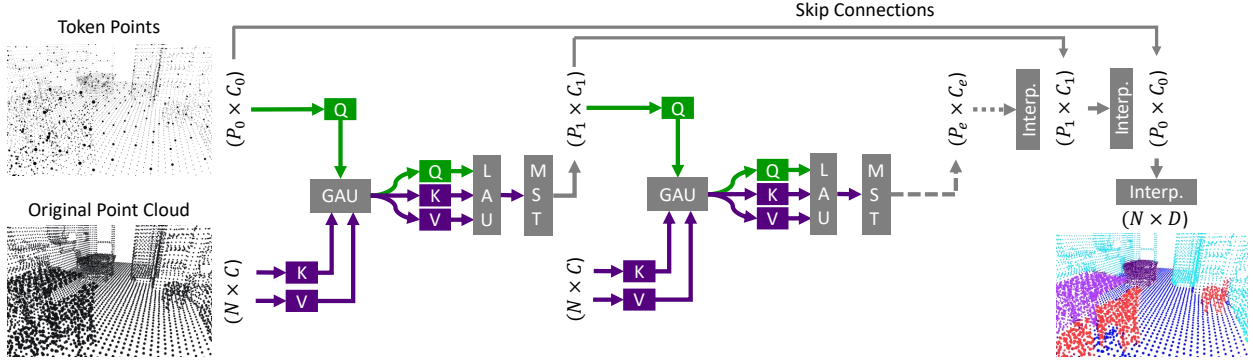


Fig. 2. Our network inputs are the original point cloud and the initial tokens (left). Local and global attention are computed to encode point clouds hierarchically (center). LAUs are involved in each local neighborhood, while GAUs are applied across the tokens and the original raw point cloud. The encoded vector is passed to interpolation layers (right) to estimate the final labels per point.

is robust to the choice of the ball radius as we tokenize a comparably large ball and rely on varying K-NNs to describe neighboring points. Contrary to previous methods that only abstract the features in a fix-sized ball, we incorporate relative distances into the output feature. Previous K-NN approaches also have an ample search space, often requiring KD-Tree search, whereas here, we do not need extra operations to find neighbors. Finally, contrary to the previous search strategies, our approach neglects irrelevant further points with the restriction imposed by the ball. We equip the representation with more flexibility by abstracting features at different scales inside the ball.

E. Network & Losses

We use the encoder model for shape classification and end-to-end segmentation tasks. The tokens produced at the last stage are pooled to a single vector for global category prediction for shape classification. For segmentation, the final stage tokens are interpolated hierarchically to recover the initial resolution of the point cloud, and the interpolated features lead to per-point predictions. We use cross-entropy loss for training both classification and segmentation models with

$$L_{cls} = -\sum_c g_c \cdot \log(p_c), \quad (4)$$

where c is the category, g_c is ground truth of the c^{th} category, and p_c is the predicted probability for the c^{th} class. We only need to average over the cross-entropy losses of all points to construct the training objective for the segmentation task:

$$L_{seg} = -\frac{1}{n} \sum_n \sum_c g_c^n \cdot \log(p_c^n), \quad (5)$$

where n is the total number of points. Ground truth is given by g_c^n for the n^{th} sample and the c^{th} class and p_c^n represents the predicted probability, respectively.

To recover the initial resolution for segmentation, we interpolate the features at each interpolation level using inverse distance weighted average [1] based on k nearest neighbors.

$$f_c = \frac{\sum_{i=1}^k w_i(c) f_i}{\sum_{i=1}^k w_i(c)} \quad \text{where} \quad w_i(c) = \frac{1}{d(c, p_i)^2} \quad (6)$$

where c is the position to be interpolated, p_i are its k nearest neighbor coordinates of Euclidean distance $d(c, p_i)$, and f_i describes the corresponding point features.

IV. DATASETS AND METRICS

This section introduces the datasets on which we performed our experiments, followed by their respective evaluation metrics.

A. ModelNet40

The ModelNet40 [22] dataset is used for shape classification tasks. ModelNet40 consists of CAD models in 40 categories. The authors split the dataset into 9843 training models and 2468 testing models. To generate point clouds from CAD models, uniform sampling is applied to each model, along with its normal vectors. We sample 1024 points from each CAD model.

The evaluation metrics adopted for the shape classification task are OA (Overall Accuracy) and mAcc (mean Accuracy). OA gives the accuracy of all test data across classes, while mAcc is based on the accuracy within each class.

B. ShapeNetPart

To conduct fine-grained 3D segmentation, we use the dataset ShapeNetPart [23]. It contains 16881 shapes, and every shape is from one of 16 classes. In total, there exist 50 unique part classes.

For evaluations, we use two metrics, category mIoU (cat. mIoU) and instance mIoU (ins. mIoU). For a shape S and its category C , we compute the intersection over union (IoU) for each part belonging to the category C based on predictions and respective ground truth. The average IoU of all part types of category C is the mIoU for shape S . By averaging mIoUs over all shapes, we also calculate ins. mIoU. Finally, for cat. mIoU, we first calculate mIoUs averaging all shapes of a given category and then average over all category mIoUs.

C. 3D Indoor Scene

We also conduct experiments on large-scale scene segmentation for real scenes. We use the Stanford 3D semantic parsing dataset [24]. It is collected by Matterport scanners from six areas, including 271 rooms. Each point belongs to one of the 13 categories. We followed the standard protocol from PointNet [2] to split training and test sets.

Each room consists of different points and is too large to fit in memory for training. Following PointNet [2], each room is split into $1\text{m} \times 1\text{m}$ blocks. During training, 4096 points are sampled from each block on the fly and are regarded as one point cloud sample. At test time, the entire scene is used for evaluation without any sampling. Three evaluation metrics are used, i.e., OA, mAcc, and mIoU. In this case, mIoU is defined as the average of IoUs over 13 classes.

V. EXPERIMENTS AND RESULTS

This section presents the results for every method, along with a comparison and analysis. Moreover, we perform extensive ablation studies on different aspects of the proposed model. Then, the latency and number of parameters are compared to justify the efficiency of our proposed methods. Finally, we illustrate the results of inferred scenes and object parts to provide qualitative evidence for our model’s performance.

We train all models on a single NVIDIA RTX 2080 TI GPU. The batch size for all tasks is set to 16. We train our network 150 epochs for the shape classification task and 200 epochs for both part segmentation and scene segmentation tasks. Cosine annealing learning rate scheduler [25] and Lamb optimizer [26] are used for optimization.

A. Shape Classification

There have been several methods benchmarking on ModelNet40 [22] shape classification. Tab. I compares our predictions with the state-of-the-art shape classification methods. As shown in the results, the proposed method achieves state-of-the-art results in terms of mAcc. The predictions of our model reach 91.4% mean accuracy, while PointTransformer [11] achieves 90.6%.

B. Part Segmentation

Tab. II shows the performance of our method for Part Segmentation on the ShapeNetPart dataset [23]. The mean IoU metric suggests the similarity of the predictions with the ground truth part semantic labels. Due to saturation on this synthetic dataset, most benchmarked works report similar performance in a small band around 0.852, for instance mean IoU. We see a performance on par with most recent works in terms of mean IoU.

C. Semantic Scene Segmentation

In this experiment, we evaluate our network performance for semantic scene segmentation. The point clouds in the Stanford 3D semantic parsing dataset [24] are from large-scale reconstructions with millions of points. Previous methods either subsample drastically, lose fine-grained details or

TABLE I
COMPARISON WITH OTHER METHODS FOR THE SHAPE CLASSIFICATION TASK ON MODELNET40 [22].

Method	mAcc	OA
3DShapeNets [22]	0.773	0.847
Perceiver [20]	-	0.857
VoxNet [12]	0.83	0.859
MVCNN [27]	-	0.901
PointNet [2]	0.862	0.892
Set Transformer [8]	-	0.904
YOGO (KNN) [9]	-	0.914
YOGO (Ball query) [9]	-	0.913
PointNet++ [1]	-	0.919
PointCNN [10]	0.881	0.922
DGCNN [4]	0.902	0.922
PointConv [3]	-	0.925
PointTransformer [11]	0.906	0.937
Ours	0.914	0.929

TABLE II
COMPARISON WITH OTHER METHODS FOR PART SEGMENTATION TASK ON SHAPENETPART DATASET [23].

Method	ins. mIoU	cat. mIoU
PointNet [2]	0.837	0.804
PCCN [28]	0.851	0.818
PointNet++ [1]	0.851	0.819
DGCNN [4]	0.851	0.823
YOGO (Ball query) [9]	0.851	-
YOGO (KNN) [9]	0.852	-
Point2Sequence [29]	0.852	-
SpiderCNN [30]	0.853	0.817
SPLATNet [31]	0.854	0.837
PointTransformer [11]	0.866	0.837
Ours	0.852	0.830

apply their networks on a significantly smaller voxel grid. We follow the inference procedure described in [9]. In Tab. III we compare with the methods that process voxels of size $1 \times 1 \times 1$. Our method can provide high-quality predictions compared to previous methods for this task and achieves 0.643 in terms of mAcc and 0.554 in mIoU. By conducting a hybrid multi-scale approach, we improve, for instance, over YOGO [9], which either uses ball query or K-NN to sample and describe points.

D. Performance Analysis

Besides the effectiveness of our proposed method on different tasks, our model can run at very low latency. Table IV details latency and memory consumption compared to previous methods. Other works’ latency and GPU memory consumption are taken directly from the evaluation of YOGO [9]. All methods are evaluated with a batch size of 8. Since YOGO runs experiments on higher-end hardware, we run our model and YOGO on a single RTX 2080TI GPU and scale the numbers accordingly. Our performance with and without the MST unit is 7.76 ms and 11.43 ms, respectively. This is significantly faster than other methods, including YOGO fastest at 21.3 ms. Improving the runtime with our full model by 46.3%, we provide faster and more accurate

TABLE III
COMPARISON BETWEEN METHODS FOR THE TASK OF SEMANTIC SCENE SEGMENTATION ON THE STANFORD 3D SEMANTIC PARSING DATASET [24].

Method	OA	mAcc	mIoU
PointNet [2]	-	0.49	0.411
SegCloud [32]	-	0.574	0.489
TangentConv [33]	-	0.622	0.526
YOGO (Ball query) [9]	-	-	0.538
YOGO (KNN) [9]	-	-	0.54
Ours	0.84	0.643	0.554

TABLE IV
EFFICIENCY COMPARISON WITH OTHER METHODS FOR PART SEGMENTATION ON SHAPEPART DATASET [23].

Method	Latency [ms]	GPU Memory [GB]
SpiderCNN [30]	170.1	6.5
PointCNN [10]	134.2	2.5
DGCNN [4]	86.7	2.4
PointNet++ [1]	77.7	2.0
RSNet [34]	73.8	0.8
PointNet [2]	21.4	1.5
YOGO (Ball) [9]	21.3	1.0
YOGO (KNN) [9]	25.6	0.9
Ours w/o MST	7.76	1.7
Ours w. MST	11.43	2.87

results than previous methods. Our memory footprint is thereby 1.7 GB and 2.87 GB for models without and with MST.

E. Ablation Studies

To understand the contribution of our pipeline components, we ablate different modules proposed in the paper. We compare the full model trained for part segmentation to models without LAU, GAU, and MST. The results are summarized in Tab. V.

1) *Effectiveness of LAU*: Ablating the LAU module, presented in section III-B, indicates that LAU improves the predictions in both cat. mIoU and ins. mIoU metrics.

2) *Effectiveness of GAU*: Similar to LAU, we ablate the effectiveness of the GAU unit described in section III-C. The results show that our model benefits from Global Attention as well.

3) *Effectiveness of MST*: Lastly, we also ablate our MST module. Compared to our full method, MST provides subtle

TABLE V
ABLATION STUDY ON THE EFFECTIVENESS OF LOCAL (LAU), GLOBAL (GAU), AND MULTI-SCALE (MST) MODULES.

LAU	GAU	MST	cat. mIoU	ins. mIoU
	✓		0.817	0.845
✓			0.816	0.847
✓	✓		0.821	0.8525
✓	✓	✓	0.830	0.8521

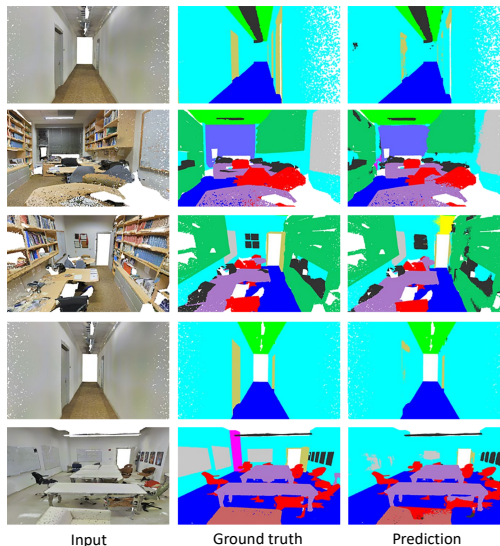


Fig. 3. Qualitative results for the scene segmentation task on Stanford 3D semantic parsing dataset [24].

improvements for part segmentation in cat. mIoU from 0.821 to 0.830. The comparison also verifies the effectiveness of our tokenization approach.

F. Qualitative Results

Besides extensive quantitative results, we also provide qualitative results and visualization of the predictions. We show sample results from our part segmentation network in Fig. 4. As illustrated, our methods predict parts classes with details very accurately.

Similar to part segmentation Fig. 3 depicts the results for scene segmentation. The input point cloud data with embedded RGB colors is shown on the right. The comparison with ground truth shows robust predictions and supports the numerical evaluations and effectiveness of the proposed multi-scale model.

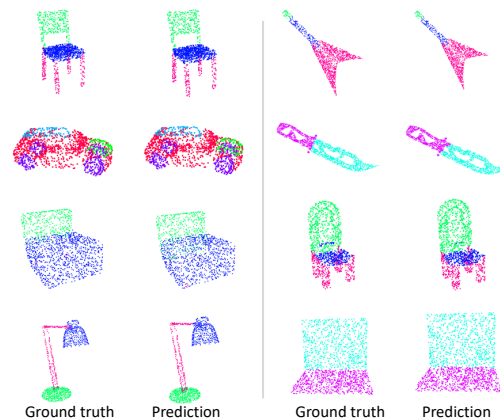


Fig. 4. Qualitative results for the semantic part segmentation task on ShapeNetPart [23] dataset.

VI. CONCLUSION

This paper presents a novel and efficient hierarchical approach to process point clouds. We design our pipeline based on attention concepts focusing on three key aspects: a) local similarity and attention help to describe local geometry with less computational burden, b) a global attention module is defined to capture global context by computing cross-attention between the sampled tokens and original cloud, and c) multi-scale tokenization describes a region by taking different levels of proximity into account. Creating a flexible and efficient point cloud processing pipeline improves standard tasks such as shape classification and semantic segmentation while improving latency. We believe that this work can pave the way to more accurate point cloud processing pipelines even on large clouds with varying densities, which is elemental for successful applications of 3D learning.

REFERENCES

- [1] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," *arXiv preprint arXiv:1706.02413*, 2017.
- [2] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [3] W. Wu, Z. Qi, and L. Fuxin, "Pointconv: Deep convolutional networks on 3d point clouds," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9621–9630.
- [4] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph cnn for learning on point clouds," *Acm Transactions On Graphics (tog)*, vol. 38, no. 5, pp. 1–12, 2019.
- [5] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [6] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [7] N. Engel, V. Belagiannis, and K. Dietmayer, "Point transformer," *IEEE Access*, vol. 9, pp. 134 826–134 840, 2021.
- [8] J. Lee, Y. Lee, J. Kim, A. Kosiorek, S. Choi, and Y. W. Teh, "Set transformer: A framework for attention-based permutation-invariant neural networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 3744–3753.
- [9] C. Xu, B. Zhai, B. Wu, T. Li, W. Zhan, P. Vajda, K. Keutzer, and M. Tomizuka, "You only group once: Efficient point-cloud processing with token representation and relation inference module," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 4589–4596.
- [10] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, "Pointcnn: Convolution on x-transformed points," *Advances in neural information processing systems*, vol. 31, pp. 820–830, 2018.
- [11] H. Zhao, L. Jiang, J. Jia, P. H. Torr, and V. Koltun, "Point transformer," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 16 259–16 268.
- [12] D. Maturana and S. Scherer, "Voxnet: A 3d convolutional neural network for real-time object recognition," in *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2015, pp. 922–928.
- [13] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4490–4499.
- [14] S. Shi, C. Guo, L. Jiang, Z. Wang, J. Shi, X. Wang, and H. Li, "Pv-rcnn: Point-voxel feature set abstraction for 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 529–10 538.
- [15] M. Saleh, S. Dehghani, B. Busam, N. Navab, and F. Tombari, "Graphite: Graph-induced feature extraction for point cloud registration," in *2020 International Conference on 3D Vision (3DV)*. IEEE, 2020, pp. 241–251.
- [16] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.
- [17] M. Saleh, S.-C. Wu, L. Cosmo, N. Navab, B. Busam, and F. Tombari, "Bending graphs: Hierarchical shape matching using gated optimal transport," *arXiv preprint arXiv:2202.01537*, 2022.
- [18] H. Yu, F. Li, M. Saleh, B. Busam, and S. Ilic, "Cofinet: Reliable coarse-to-fine correspondences for robust pointcloud registration," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [19] P. Ruhkamp, D. Gao, H. Chen, N. Navab, and B. Busam, "Attention meets geometry: Geometry guided spatial-temporal attention for consistent self-supervised monocular depth estimation," in *2021 International Conference on 3D Vision (3DV)*. IEEE, 2021, pp. 837–847.
- [20] A. Jaegle, F. Gimeno, A. Brock, A. Zisserman, O. Vinyals, and J. Carreira, "Perceiver: General perception with iterative attention," *arXiv preprint arXiv:2103.03206*, 2021.
- [21] M.-H. Guo, J.-X. Cai, Z.-N. Liu, T.-J. Mu, R. R. Martin, and S.-M. Hu, "Pct: Point cloud transformer," *arXiv preprint arXiv:2012.09688*, 2020.
- [22] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3d shapenets: A deep representation for volumetric shapes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1912–1920.
- [23] L. Yi, V. G. Kim, D. Ceylan, I.-C. Shen, M. Yan, H. Su, C. Lu, Q. Huang, A. Sheffer, and L. Guibas, "A scalable active framework for region annotation in 3d shape collections," *ACM Transactions on Graphics (ToG)*, vol. 35, no. 6, pp. 1–12, 2016.
- [24] I. Armeni, O. Sener, A. R. Zamir, H. Jiang, I. Brilakis, M. Fischer, and S. Savarese, "3d semantic parsing of large-scale indoor spaces," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1534–1543.
- [25] I. Loshchilov and F. Hutter, "Sgdr: Stochastic gradient descent with warm restarts," *arXiv preprint arXiv:1608.03983*, 2016.
- [26] Y. You, J. Li, S. Reddi, J. Hseu, S. Kumar, S. Bhojanapalli, X. Song, J. Demmel, K. Keutzer, and C.-J. Hsieh, "Large batch optimization for deep learning: Training bert in 76 minutes," *arXiv preprint arXiv:1904.00962*, 2019.
- [27] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, "Multi-view convolutional neural networks for 3d shape recognition," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 945–953.
- [28] S. Wang, S. Suo, W.-C. Ma, A. Pokrovsky, and R. Urtasun, "Deep parametric continuous convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2589–2597.
- [29] X. Liu, Z. Han, Y.-S. Liu, and M. Zwicker, "Point2sequence: Learning the shape representation of 3d point clouds with an attention-based sequence to sequence network," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 8778–8785.
- [30] Y. Xu, T. Fan, M. Xu, L. Zeng, and Y. Qiao, "Spidernn: Deep learning on point sets with parameterized convolutional filters," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 87–102.
- [31] H. Su, V. Jampani, D. Sun, S. Maji, E. Kalogerakis, M.-H. Yang, and J. Kautz, "Splatnet: Sparse lattice networks for point cloud processing," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2530–2539.
- [32] L. Tchapmi, C. Choy, I. Armeni, J. Gwak, and S. Savarese, "Segcloud: Semantic segmentation of 3d point clouds," in *2017 international conference on 3D vision (3DV)*. IEEE, 2017, pp. 537–547.
- [33] M. Tatarchenko, J. Park, V. Koltun, and Q.-Y. Zhou, "Tangent convolutions for dense prediction in 3d," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3887–3896.
- [34] Q. Huang, W. Wang, and U. Neumann, "Recurrent slice networks for 3d segmentation of point clouds," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2626–2635.