

Hierarchical Point Cloud Encoding and Decoding with Lightweight Self-Attention based Model

En Yen Puang¹, Hao Zhang², Hongyuan Zhu¹ *Member, IEEE*, and Wei Jing^{1,3}

Abstract—In this paper we present SA-CNN, a hierarchical and lightweight self-attention based encoding and decoding architecture for representation learning of point cloud data. The proposed SA-CNN introduces convolution and transposed convolution stacks to capture and generate contextual information among unordered 3D points. Following conventional hierarchical pipeline, the encoding process extracts feature in local-to-global manner, while the decoding process generates feature and point cloud in coarse-to-fine, multi-resolution stages. We demonstrate that SA-CNN is capable of a wide range of applications, namely classification, part segmentation, reconstruction, shape retrieval, and unsupervised classification. While achieving the state-of-the-art or comparable performance in the benchmarks, SA-CNN maintains its model complexity several order of magnitude lower than the others. In term of qualitative results, we visualize the multi-stage point cloud reconstructions and latent walks on rigid objects as well as deformable non-rigid human and robot models.

Index Terms—Visual Learning, Deep Learning for Visual Perception, Recognition

I. INTRODUCTION

WITH recent advancement on 3D sensory devices and the surging needs on 3D perception in applications such as the robotics and autonomous vehicles [1], understanding point cloud data from 3D sensors becomes increasingly important. Unlike 2D pixel-based images, 3D point cloud data is unordered and unstructured. Moreover, point clouds are usually in large scale and thus effective representation learning is often desirable for many downstream tasks [2].

Many deep learning based approaches [3] have been applied to 3D point cloud perception. Early work rely on handcrafted 3D features which require intensive feature engineering and are difficult to generalize. Then methods like [4], [5] convert 3D point cloud into voxel representation and extend 2D convolution to 3D, but bear high computational cost and low voxel resolution. Some work [6], [7] explore multi-view projection methods that apply 2D CNN to extract the view-wise features circumvent the heavy computational complexity, but

Manuscript Received September 9 2021; Revised December 30 2021; Accepted January 25 2022.

This paper was recommended for publication by Editor Cesar C. Lerma upon evaluation of the Associate Editor and Reviewers' comments. This work was supported by AME Programmatic Funding Scheme (Project #A18A2b0046), Career Development Fund (Project C210812033) and RobotHTPO Seed Fund (Project C211518008).

¹E. Y. Puang, H. Zhu and W. Jing are with Institute for Infocomm Research, A*STAR, Singapore {puang_en_yen, zhuh}@i2r.a-star.edu.sg

²H. Zhang is with Centre for Frontier AI Research, A*STAR, Singapore zhang_hao@ihpc.a-star.edu.sg

³W. Jing is also with Alibaba Group, China 21wjing@gmail.com
Digital Object Identifier (DOI): 10.1109/LRA.2022.3149569

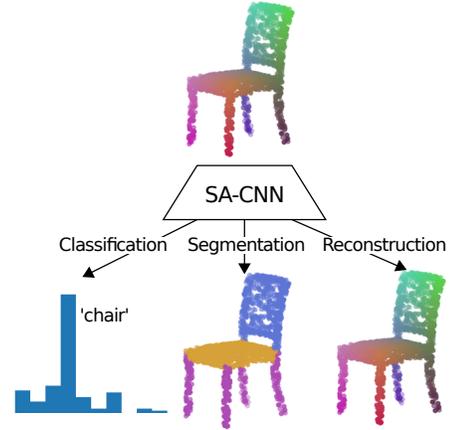


Fig. 1. SA-CNN is a lightweight self-attention based architecture capable for point cloud analysis tasks such as classification, segmentation and reconstruction by auto-encoding.

usually requires heuristics in choosing informative viewpoints. Later, [8], [9] propose an efficient and effective end-to-end architecture that directly learns the representations from the unstructured point cloud data.

In term of dimensionality reduction, auto-encoder or encoder-decoder architecture has been widely adopted in compressing and/or distilling input into lower dimensional representation in multiple domains. In the context of point cloud, earlier methods [8], [10], [2] use mostly MLP in feature extraction and generation. Hierarchical pipeline was later introduced [9], [11], [12] to encode/extract feature in local-to-global manner, and to decode/generate point cloud in coarse-to-fine, multi-resolution stages.

To this end, we propose to apply a lightweight hierarchical self-attention mechanism [13], [9] to handle the unstructured, interrelated contextual information of point cloud data. Self-attention models have been widely adopted in natural language processing, vision and point cloud [14], [15] due to their efficiency in handling dependencies in sequential/unordered data. Leveraging on the efficiency and simplicity of the self-attention mechanism, we propose a lightweight point cloud convolution and transposed convolution operator for hierarchical feature extraction and generation on point cloud data. We demonstrate that our models are capable for a wide range of applications and achieve the state-of-the-art or comparable performance while keeping the model complexity low. Hence, the main contributions of this work are summarized as follows:

- we propose a hierarchical network architecture named SA-CNN for point cloud encoding and decoding using

an lightweight self-attention based model;

- we demonstrate the applications of SA-CNN in classification, part segmentation, reconstruction, shape retrieval and unsupervised classification;
- we show that SA-CNN achieves the state-of-the-art or comparable performances in multiple benchmarks while keeping the model complexity orders of magnitude lower.

II. RELATED WORK

A. Point-based Methods for Point Cloud Analysis

Point-based methods became popular due to its efficiency, flexibility and scalability. Different network architectures e.g., Multi-Layer Perceptron (MLP), point convolution, graph-based, and attention-based network could be used to deal with the point-wise feature extraction. [8], [9] propose PointNet(++) as pioneer work for point-based methods on 3D point cloud data, which directly learn feature representation from the point cloud data. Convolution-based methods have also been applied to point cloud data [16], [17] by directly performing convolution-like operation on 3D points. More recently, graph-based neural networks [18], [19] and attention-based methods [20], [21], [22] also demonstrate good performance for point cloud feature learning. However, these related works only learn features in encoding manner. Our work is a point-based method that uses self-attention mechanism to build hierarchical encoding and decoding models for efficient representation learning of point cloud.

B. Self-Attention Mechanism

Self-attention mechanism has become an integral part of various neural network-based models for handling the sequence inputs, allowing modeling of dependencies without regard to their distance in the sequences. Specifically, the Transformers [13] have revolutionized natural language processing tasks due to their tremendous ability on representation learning. Recently, there are several work [14], [15] adopt the attention mechanism to address the 3D point cloud problems. [23], [24], [25] proposed additional mechanism together with self-attention to improve point sampling and feature aggregation. These self-attention based methods have achieved promising performance on various point cloud tasks, but their models are still less efficient with a large amount of parameters and high complexity. Our methods adopt the same self-attention mechanism in the design of a lightweight architecture to achieve the state-of-the-art or comparable performance with orders of magnitude lesser model parameter and computational complexity.

C. Point Cloud Auto-Encoding

Point cloud generation are generally divided into 2 categories: Generative Adversarial Network (GAN) based and non-GAN based methods. GAN-based generative model [2], [26], [12], [27] sample point cloud based on the data distribution learned during training. These methods often output realistic point cloud with great details, but generally require more complex training procedures.

Non-GAN based methods are mostly in encoder-decoder, unsupervised architecture [28]. To improve reconstruction performance, [29] requires the point cloud to be a 1D ordered list structure. [11], [10] parameterize the decoder with 2D grid patches as part of the latent representation, while [27] parameterizes point cloud with spectral frequencies. Our auto-encoder is a non-GAN, self-supervised based architecture that reconstructs input point cloud based only on a latent vector and without any other parameterization.

III. PRELIMINARY

A. Self-Attention Layer

The self-attention layer $\mathcal{A}(\cdot)$ we adopt is a scaled dot-product attention mechanism [13] $\mathcal{A} : Q, K, V \rightarrow \mathbb{R}^{N \times d'}$ which takes as input *Query* $Q \in \mathbb{R}^{N \times d}$, *Key* $K \in \mathbb{R}^{N' \times d}$ and *Value* $V \in \mathbb{R}^{N' \times d'}$, and aggregates information from V for each token in Q based on the alignment with K :

$$\mathcal{A}(Q, K, V) = \sigma \left(\frac{QK^\top}{\sqrt{d}} \right) V \quad (1)$$

where σ is a softmax operator. Multi-head attention simply concatenates the output of multiple self-attention layer. The output is of $\mathbb{R}^{N \times hd'}$ where h is the number of head.

Point cloud data is here defined as an unordered set of N points $\mathcal{P} = \{\mathbf{p}_i\}_{i=1}^N \mid \mathbf{p} \in \mathbb{R}^3$ in 3D Euclidean space without additional features. When applied on point cloud processing, self-attention layer acts analogously to a convolution kernel. *Query* points are analogous to the kernel's center. Whereas the neighbouring points, in which the *Key* and *Value* are derived from, are analogous to the rest of the kernel inputs. Multi-head attention is then analogous to a convolution filter. Self-attention model is well suited in handling point cloud [14], [15], [30], [31] because:

- Self-attention is invariant to input permutation and hence able to handle input point cloud data in the form of unordered set.
- Instead of one-to-many relationship rendered by convolution kernel, self-attention improves points interactions by having many-to-many relationship in the input set in an efficient manner.
- With adequate normalization in the input set, self-attention is robust under rigid transformation such as translation and rotation.

B. Processing Pipeline

Unlike CNN's grid structure, point cloud requires several pre-processing layers in a hierarchical pipeline [9]:

- *Sampling*. Given a set of points, a subset is selected and formed into an unstructured probing grid.
- *Grouping*. Given a set of points and its probing grid, the nearest neighbours are gathered into local neighbourhoods for each probing points.
- *Normalizing*. Given a set of local neighbourhoods, points are normalized by subtracting its probing point and scaling by a constant.
- *Interpolating*. Given neighbourhoods from 2 adjacent levels, features are propagated from the smaller neighbourhood with distance based interpolation.

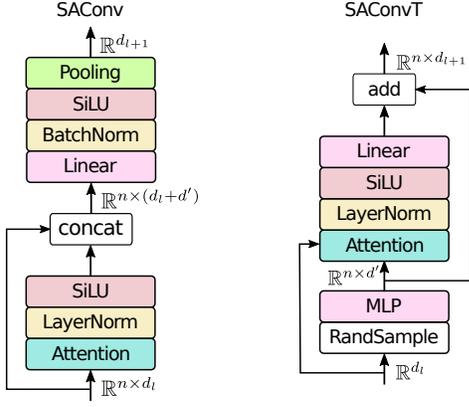


Fig. 2. The design of SACConv and SACConvT stacks. The former aggregates features from a set, while the latter expands a feature into a set.

IV. MODEL ARCHITECTURE

The proposed SA-CNN consists of 2 basic operators for processing data in the form of unordered set. We will first discuss these operators, then we describe how to build with these operators for different applications.

A. SA-CNN Modules

SACConv is a set convolution stack that aggregates features $\text{SACConv} : \mathbb{R}^{n \times d_i} \rightarrow \mathbb{R}^{d_{i+1}}$ from an unordered set of size n and outputs a single feature vector. For a self-attention layer to work analogous to a convolution kernel on point cloud, *Sampling-Grouping-Normalizing* (SGN) pre-processing layer is deployed to form the input $\mathcal{X} \in \mathbb{R}^{(k+1) \times d}$ where d is the feature dimension and k is the number of nearest neighbour of a probing point. *Query*, *Key* and *Value* in Eq 1 are from

$$Q = \mathcal{X}W_q \quad (2a)$$

$$K = \mathcal{X}W_k \quad (2b)$$

$$V = \mathcal{X}W_v \quad (2c)$$

the linear layers where $W_{q,k} \in \mathbb{R}^{d \times d_q}$ and $W_v \in \mathbb{R}^{d \times d'}$ are the weights. The output of the multi-head attention then goes through a sequence of common layers as depicted in Fig. 2 (left). A pooling layer is placed at the end as the symmetric function that aggregates the final features in the set.

SACConvT is a transposed convolution stack proposed to generate an unordered set $\text{SACConvT} : \mathbb{R}^{d_i} \rightarrow \mathbb{R}^{n \times d_{i+1}}$ based on a single input feature vector \mathcal{X} . It consists of a random sampling layer that draw $\mathcal{N} \in \mathbb{R}^{n \times d'}$ samples from unit Gaussian distribution and

$$\mathcal{X}' = \text{MLP}\left(\text{concat}_d(\mathcal{X}, \mathcal{X}W_r * \mathcal{N})\right) \quad (3)$$

where $W_r \in \mathbb{R}^{d \times d'}$ is the weight of the linear layer before the broadcasted element-wise multiplication $*$ with \mathcal{N} , and concat_d is concatenation along feature dimension. Multi-head attention takes it as input $\mathcal{A}(\mathcal{X}')$ and then followed by a sequence of common layers as depicted in Fig 2 (right).

Positional Encoding were proposed [13] to provide sequential information for each tokens. Although unordered, points in point cloud are well represented in euclidean space. Therefore,

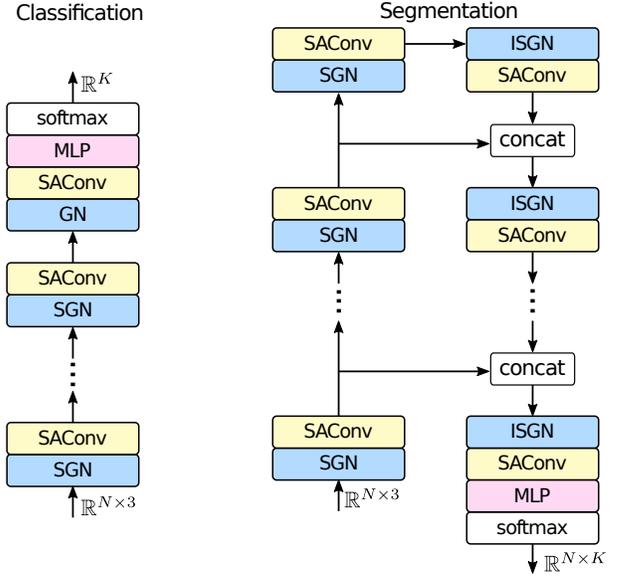


Fig. 3. The architectures of classification and segmentation networks. Both networks use SACConv paired with pre-processing layers [9] (I/SG/N) in each hierarchical levels.

we use the 3D xyz coordinate of point as its positional encoding. Combining *Normalizing* pre-processing layer, the positional encoding is its normalized coordinate with respect to the probing point of the neighbourhood. In the hierarchical pipeline, $d_l = 3 + d_{l-1}$ given $d_0 = 0$ indicates the input dimension of level l . Concretely, in each hierarchical levels the input to SACConv is the concatenation of the normalized 3D coordinate and the features computed from the previous level given the input dimension is 3.

B. Applications of SA-CNN

1) *Classification and Part Segmentation*: We apply SA-CNN modules as a drop-in replacement to solve the point cloud classification and part segmentation tasks using the hierarchical feature extraction pipeline [9] as depicted in Fig 3. Classification network is composed of several levels, each consists of a SGN-SACConv pair. Each level abstracts a subset of probing points, while the last level groups (GN, without *sampling*) all the remaining for the final classification output. The part segmentation network adopts U-Net [32] structure with skip connections. ISGN combines *interpolating* layer with SGN and is used to propagate features to the next level.

2) *Auto-encoding*: We apply SA-CNN modules to solve the point cloud reconstruction task with a hierarchical auto-encoder architecture depicted in Fig 4. Similar to the pipeline of the classification task, the **encoder** $f : \mathcal{P} \mapsto z$ takes in as input a set of points $\mathcal{P} \in \mathbb{R}^{N \times 3}$ and processes it with several SGN-SACConv pairs. A MLP then turns the features into a latent representation vector $z \in \mathbb{R}^L$ where $L \ll N$.

The hierarchical **decoder** $g : z \mapsto \{\bar{\mathcal{P}}_l\}_{l=1}^H$ takes as input the latent vector, and expands it through H levels of SG-SACConvT pairs. Each hierarchical level expands every points in the level and produces a set of output points $\bar{\mathcal{P}}$ given $|\bar{\mathcal{P}}_a| < |\bar{\mathcal{P}}_b|$ if $a < b$ and $|\bar{\mathcal{P}}_H| = N$.

TABLE I
MODELNET40 CLASSIFICATION BENCHMARK. INPUT CONSISTS OF
NUMBER OF POINTS, 'P'(POINT) AND/OR 'N'(NORMAL).

Methods	Input	#param(M)	Accuracy(%)
PointNet [8]	1k p	3.5	89.2
PointNet++ [9]	1k p	1.5	90.7
PointNet++ [9]	5k p,n	1.5	91.9
Kd-Net [35]	32k p	2.0	91.8
PointCNN [16]	1k p	0.3	92.5
DGCNN [18]	1k p	1.8	92.9
PCNN [36]	1k p	1.4	92.3
DensePoint[33]	1k p	0.7	93.2
RS-CNN[37]	1k p	1.4	93.6
KPConv [17]	7k p	14.3	92.9
Pnt Tfmer [24]	1k p,n	13.50	92.8
Pnt Tfmer [20]	1k p	9.6	93.7
PCT [21]	1k p	2.9	93.2
PointASNL [25]	1k p	-	92.9
AttPNet [22]	1k p	-	93.6
PAT [23]	1k p	-	91.7
SA-CNN (Ours)	1k p	0.04	92.3

TABLE II
COMPUTATIONAL COMPLEXITY FOR MODELNET40 CLASSIFICATION. ALL
METHODS TAKE AS INPUT 1024 POINTS ON A GEFORCE RTX 3090.

Methods	FLOPs /sample	Params size(Mb)	Training size(Mb)
PointNet [8]	878M	14	60
PointNet++ [9]	1.69G	6	138
DGCNN [18]	4.78G	7	185
KPConv [17]	200M	-	-
PointCNN [16]	210M	1	22
Pnt Tfmer [24]	92G	88	5736
Pnt Tfmer [20]	3.6G	38	6334
PCT [21]	70G	12	2400
DensePoint [33]	651M	-	-
RS-CNN [37]	295M	-	-
SA-CNN (ours)	8M	0.2	2

axes, scaling between $[0.6667, 1.5]$ and translation between $[-0.2, 0.2]$ in all 3 axes. During testing, we adopt voting strategy where 10 tests are performed with the same data augmentations except rotations and noise.

Table I shows the comparison of classification results in term of overall accuracy and the number of parameters on ModelNet40. Statistic in term of floating point operations (FLOPs) per sample and memory for parameter and training required on GPU are depicted in Table II. SA-CNN classifier has 4/5/6/7 as number of head and attention size, 8/7/6/16 as the neighbourhood size, with sampling ratio $0.3N/0.075N/0.015N$, feature size 16/32/64/96 in the 4 levels, and a $p = 0.6$ dropout before the output layers.

Due to the simple and efficient hierarchical architecture SA-CNN achieves the lowest number of parameter, FLOPs and memory by several order of magnitude, while having comparable accuracy to the other self-attention and non-attention methods. This advantage makes SA-CNN suitable for low computation/memory point cloud applications. Although

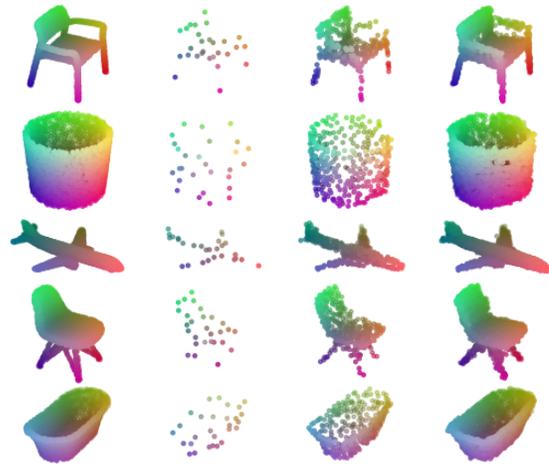


Fig. 6. Hierarchical point cloud decoding on ShapeNet with 2048 points. Left column is the ground truth and right columns are the reconstructions.



Fig. 7. Hierarchical point cloud decoding on robots and human. Left column is the ground truth and right columns are the reconstructions.

generic, the accuracy gap (1.4%) to the state-of-the-art suggests that SA-CNN could work with additional mechanism [25], [24], [23] to boost classification performance.

B. Part Segmentation

We evaluate SA-CNN on fine-grained object recognition using ShapeNet part [39] benchmark in point cloud parts segmentation task which consists of 12137 models for training and 2874 for testing. Models are classified into 16 categories, and each category contains between 2 to 6 parts. For each model we zero the bounding box's center and normalize to a unit sphere. During training we augment data with random rotation between $[-10, 10]$ degree, scaling between $[0.8, 1.25]$ and translating between $[-0.1, 0.1]$ in all 3 axes. The performance metric used for this task is intersection-over-union (IoU). There are 3 IoUs reported: IoU for each categories, mean IoU (mIoU) over all categories, and mIoU over all instances.

TABLE III
PART SEGMENTATION RESULTS (IoU %) ON SHAPENET PART BENCHMARK. INPUTS ARE EITHER 3D POINTS, OR POINTS WITH NORMAL (N).

method	input	cat. mIoU	int. mIoU	air plane	bag	cap	car	chair	ear phone	guitar	knife	lamp	laptop	motor bike	mug	pistol	rocket	skate board	table
SCN[38]	1k	81.8	84.6	83.8	80.8	83.5	79.3	90.5	69.8	91.7	86.5	82.9	96.0	69.2	93.8	82.5	62.9	74.4	80.8
Kd-Net[35]	4k	77.4	82.3	80.1	74.6	74.3	70.3	88.6	73.5	90.2	87.2	81.0	94.9	57.4	86.7	78.1	51.8	69.9	80.3
PointNet[8]	2k	80.4	83.7	83.4	78.7	82.5	74.9	89.6	73.0	91.5	85.9	80.8	95.3	65.2	93.0	81.2	57.9	72.8	80.6
PointNet++[9]	2k,n	81.9	85.1	82.4	79.0	87.7	77.3	90.8	71.8	91.0	85.9	83.7	95.3	71.6	94.1	81.3	58.7	76.4	82.6
PCNN[36]	2k	81.8	85.1	82.4	80.1	85.5	79.5	90.8	73.2	91.3	86.0	85.0	95.7	73.2	94.8	83.3	51.0	75.0	81.8
DGCNN[18]	2k	82.3	85.1	84.2	83.7	84.4	77.1	90.9	78.5	91.5	87.3	82.9	96.0	67.8	93.3	82.6	59.7	75.5	82.0
RS-CNN[37]	2k	84.0	86.2	83.5	84.8	88.8	79.6	91.2	81.1	91.6	88.4	86.0	96.0	73.7	94.1	83.4	60.5	77.7	83.6
DensePnt[33]	2k	84.2	86.4	84.0	85.4	90.0	79.2	91.1	81.6	91.5	87.5	84.7	95.9	74.3	94.6	82.9	64.6	76.8	83.7
Point Tformer[20]	-	83.7	86.6	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
AttPNet [22]	2k	82.8	85.2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
GAPNet [14]	2k	82.0	84.7	84.2	84.1	88.8	78.1	90.7	70.1	91.0	87.3	83.1	96.2	65.9	95.0	81.7	60.7	74.9	80.8
LAE-Conv [15]	2k	84.1	85.9	83.3	86.1	85.7	80.3	90.5	82.7	91.5	88.1	85.5	95.9	77.9	95.1	84.0	64.3	77.6	82.8
SA-CNN (Ours)	2k	84.0	86.7	85.1	85.3	85.4	78.6	91.7	79.4	91.4	87.6	86.8	95.5	72.2	94.6	81.3	62.3	81.5	84.8

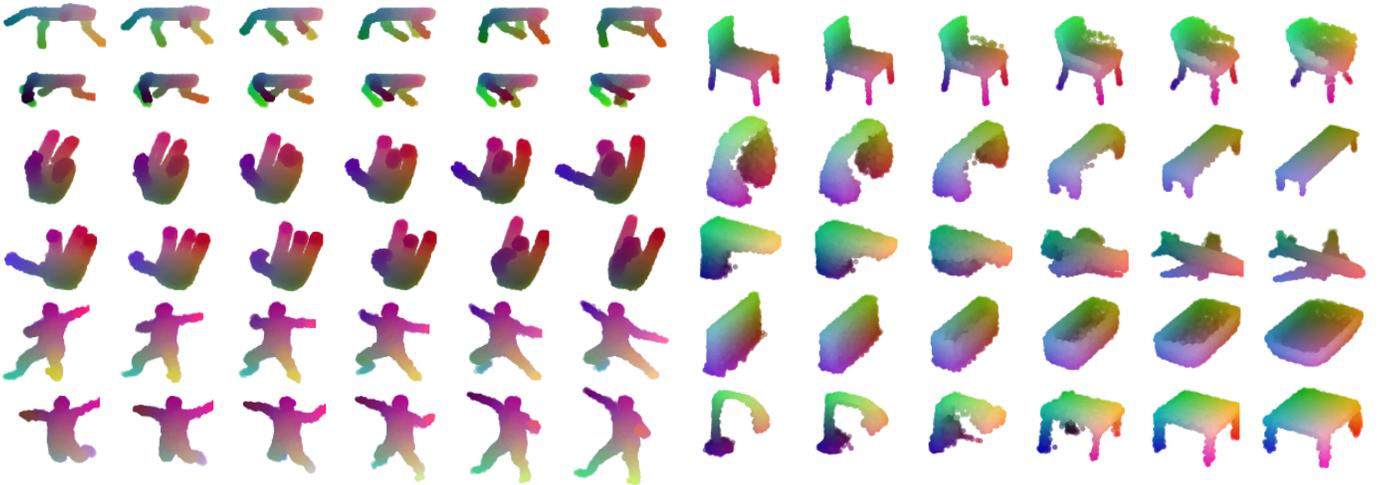


Fig. 8. Latent walk in robots and human models.

Quantitative comparison among the other methods are summarised in Table III. SA-CNN achieves state-of-the-art on instance mIoU and several category IoUs. These results demonstrate SA-CNN’s capability and robustness in recognising a diverse range of fine-grained shapes. Fig 10 depicts some examples of the ShapeNet segmentation results.

C. Reconstruction with Auto-Encoder

In this section, we demonstrate the performance of SA-CNN in auto-encoding. We train the auto-encoder using models from two different domains. For rigid object reconstruction, we use models from ModelNet40 [34] and ShapeNetCore [40] with 40 and 55 classes respectively. For non-rigid object reconstruction, we use SMPL [41] human, allegro hand and aliengo quadruped robot model. All non-rigid dataset consists of 4096 synthetically generated training models of different poses. We uniformly sample 2048 points as input as well as output for the training of auto-encoder, and visualize the hierarchical decoding and reconstructions with latent size of 128 in Fig 6 and Fig 7.

Fig. 9. Latent walk among ShapeNetCore objects.

TABLE IV
THE SHAPE RETRIEVAL RESULTS ON THE MODELNET40.

Methods	Latent Size	mAP (%)
PointNet[8]	-	70.5
PointCNN[16]	-	83.8
DGCNN[18]	-	85.3
DensePnt[33]	256	<u>88.5</u>
AE	32	84.0
AE	64	85.4
AE	128	86.7
AE	256	87.1
AE + Triplet	32	<u>88.9</u>
AE + Triplet	64	89.2
AE + Triplet	128	90.1
AE + Triplet	256	89.7

Next, we evaluate SA-CNN qualitatively on auto-encoding by visualizing the latent walks of both rigid and non-rigid objects. Linear interpolations of latent code between 2 models are fed to the decoder to get the output reconstructions. Fig 8 and Fig 9 depict some examples of latent walks.

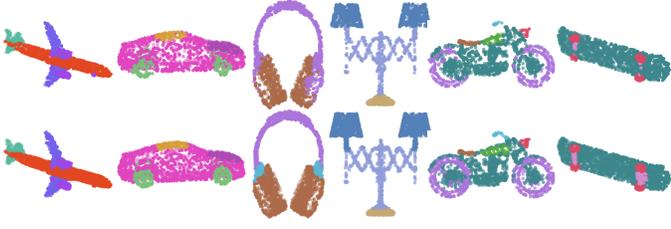


Fig. 10. Examples of ShapeNet part segmentation which contains 16 categories, each having 2 to 6 parts. Top row is segmentation predictions by SA-CNN model, and the bottom row is the corresponding ground truths.



Fig. 11. Examples of ModelNet Shape retrieval task. Given a query model (left) from the test set, top nearest neighbours (right) from the training set are retrieved. Neighbours from the wrong class are highlighted in red.

D. Shape Retrieval

We evaluate the performance of SA-CNN auto-encoding with shape retrieval task. We train the auto-encoder using ModelNet40. For this task, given a query shape in test set, we first extract its latent vector using the encoder. Then we search for the top 10 nearest latent vectors seen during training by cosine distance. We report the mean average precision (mAP) based on the retrieved nearest neighbours.

Table IV summarizes the shape retrieval performance trained with various latent sizes. With pure unsupervised training, our auto-encoder is able to match with methods that train with classification model. With just 32 latent size our method that train with reconstruction and triplet loss surpasses all other methods. Fig 11 depicts some examples of queries and the corresponding nearest neighbours. It is observed that some instances from classes such as cup and pot, and chair and stool, are highly similar, which leads to lower performance in classification by retrieval.

E. Unsupervised Classification

In this section, we evaluate the quality of SA-CNN representation learning by comparing classification accuracy using only the latent code as the input feature of a linear Support Vector Machine (SVM). We pre-train the auto-encoder using ShapeNetCore, and then fit and test the SVM using the latent

TABLE V
UNSUPERVISED CLASSIFICATION ON MODELNET40 USING LATENT REPRESENTATION PRE-TRAINED ON SHAPENETCORE.

Methods	Latent Size	Accuracy (%)
VConv-DAE [42]	6912	75.5
3D-GAN [43]	448	83.3
MRTNet-VAE [29]	224	86.4
AE-CD [2]	512	84.5
FoldingNet [11]	512	88.4
SA-CNN (Ours)	128	87.1
SA-CNN (Ours)	256	88.6

TABLE VI
SA-CNN'S PART SEGMENTATION AND AUTO-ENCODING MODEL COMPLEXITY IN NUMBER OF PARAMETER AND FLOPS PER SAMPLE.

Network	Setting	#param	#FLOPs
Part Segmentation	1k points	227k	1.2G
	2k points	227k	2.4G
Encoder/Decoder	32 latent	22k/35k	15M/20M
	64 latent	58k/77k	38M/80M
	128 latent	160k/190k	100M/180M
	128 latent 2k points	160k/190k	116M/271M

code of ModelNet40 train/test split. Table V depicts the latent size and results of the task. Our auto-encoder achieves the state-of-the-art with only 256 latent size show that SA-CNN is capable of compact and efficient unsupervised representation learning for point cloud data.

F. Model Complexity

Model complexity is an important aspect while using self-attention mechanism. In addition to Table II, we further summaries model complexity in term of number of parameters and FLOPs per sample of SA-CNN's part segmentation and auto-encoding models in Table VI.

SA-CNN part segmentation network has 4/5/6/7/8/8/7/6 number of head and attention size, $0.4N$, $0.12N$, $0.024N$ sampling ratio, 14/12/10/8 neighbourhood size, 16/32/64/96/128/128/96/64 feature dimension in the levels, and uses Max pooling for the first 3 and Avg pooling for the rest of the SAConv. Auto-encoding networks have similar settings but increase proportionally to the latent size.

VI. CONCLUSION

In this work, we present SA-CNN, a lightweight self-attention based encoding and decoding architecture for 3D point cloud representation learning. The proposed SA-CNN is a drop-in replacement for multiple hierarchical point cloud analysis tasks. We demonstrate that SA-CNN achieves the state-of-the-art or comparable performance to other methods on multiple benchmarks with significantly lower model complexity. In addition, we visualize the multi-stages point cloud reconstruction and their latent walks on rigid objects as well as non-rigid human and robot models.

REFERENCES

- [1] E. Potapova, M. Zillich, and M. Vincze, "Survey of recent advances in 3d visual attention for robotics," *The International Journal of Robotics Research*, vol. 36, no. 11, pp. 1159–1176, 2017.
- [2] P. Achlioptas, O. Diamanti, I. Mitliagkas, and L. Guibas, "Learning representations and generative models for 3d point clouds," in *International conference on machine learning*. PMLR, 2018, pp. 40–49.
- [3] Y. Guo, H. Wang, Q. Hu, H. Liu, L. Liu, and M. Bennamoun, "Deep learning for 3d point clouds: A survey," *IEEE transactions on pattern analysis and machine intelligence*, 2020.
- [4] 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.
- [5] G. Riegler, A. Osman Ulusoy, and A. Geiger, "Octnet: Learning deep 3d representations at high resolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3577–3586.
- [6] 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.
- [7] Z. Yang and L. Wang, "Learning relationships for multi-view 3d object recognition," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 7505–7514.
- [8] 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.
- [9] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," *Advances in Neural Information Processing Systems*, vol. 30, p. 5099–5108, 2017.
- [10] T. Groueix, M. Fisher, V. G. Kim, B. C. Russell, and M. Aubry, "A papier-mâché approach to learning 3d surface generation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 216–224.
- [11] Y. Yang, C. Feng, Y. Shen, and D. Tian, "Foldingnet: Point cloud auto-encoder via deep grid deformation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 206–215.
- [12] D. W. Shu, S. W. Park, and J. Kwon, "3d point cloud generative adversarial network based on tree structured graph convolutions," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 3859–3868.
- [13] 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.
- [14] C. Chen, L. Z. Fragonara, and A. Tsourdos, "Gapointnet: Graph attention based point neural network for exploiting local feature of point cloud," *Neurocomputing*, vol. 438, pp. 122–132, 2021.
- [15] M. Feng, L. Zhang, X. Lin, S. Z. Gilani, and A. Mian, "Point attention network for semantic segmentation of 3d point clouds," *Pattern Recognition*, vol. 107, p. 107446, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0031320320302491>
- [16] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, "Pointcnn: Convolution on χ -transformed points," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2018, pp. 828–838.
- [17] H. Thomas, C. R. Qi, J.-E. Deschaud, B. Marcotegui, F. Goulette, and L. J. Guibas, "Kpconv: Flexible and deformable convolution for point clouds," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6411–6420.
- [18] 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.
- [19] G. Te, W. Hu, A. Zheng, and Z. Guo, "Rgcnn: Regularized graph cnn for point cloud segmentation," in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 746–754.
- [20] 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.
- [21] M.-H. Guo, J.-X. Cai, Z.-N. Liu, T.-J. Mu, R. R. Martin, and S.-M. Hu, "Pct: Point cloud transformer," *Computational Visual Media*, vol. 7, no. 2, pp. 187–199, 2021.
- [22] Y. Yang, Y. Ma, J. Zhang, X. Gao, and M. Xu, "Attnpnet: Attention-based deep neural network for 3d point set analysis," *Sensors*, vol. 20, no. 19, p. 5455, 2020.
- [23] J. Yang, Q. Zhang, B. Ni, L. Li, J. Liu, M. Zhou, and Q. Tian, "Modeling point clouds with self-attention and gumbel subset sampling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3323–3332.
- [24] N. Engel, V. Belagiannis, and K. Dietmayer, "Point transformer," *IEEE Access*, vol. 9, pp. 134 826–134 840, 2021.
- [25] X. Yan, C. Zheng, Z. Li, S. Wang, and S. Cui, "Pointasnl: Robust point clouds processing using nonlocal neural networks with adaptive sampling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5589–5598.
- [26] C.-L. Li, M. Zaheer, Y. Zhang, B. Póczos, and R. Salakhutdinov, "Point cloud gan," *arXiv preprint arXiv:1810.05795*, 2018.
- [27] S. Ramasinghe, S. Khan, N. Barnes, and S. Gould, "Spectral-gans for high-resolution 3d point-cloud generation," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 8169–8176.
- [28] H. Fan, H. Su, and L. J. Guibas, "A point set generation network for 3d object reconstruction from a single image," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 605–613.
- [29] M. Gadelha, R. Wang, and S. Maji, "Multiresolution tree networks for 3d point cloud processing," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 103–118.
- [30] G. Zhang, Q. Ma, L. Jiao, F. Liu, and Q. Sun, "Attn: attention adversarial networks for 3d point cloud semantic segmentation," in *Proceedings of the Twenty-Ninth International Conference on Artificial Intelligence*, 2021, pp. 789–796.
- [31] X. Wang, J. He, and L. Ma, "Exploiting local and global structure for point cloud semantic segmentation with contextual point representations," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [32] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [33] Y. Liu, B. Fan, G. Meng, J. Lu, S. Xiang, and C. Pan, "Densepoint: Learning densely contextual representation for efficient point cloud processing," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5239–5248.
- [34] 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.
- [35] R. Klokov and V. Lempitsky, "Escape from cells: Deep kd-networks for the recognition of 3d point cloud models," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 863–872.
- [36] M. Atzmon, H. Maron, and Y. Lipman, "Point convolutional neural networks by extension operators," *ACM Transactions on Graphics (TOG)*, vol. 37, no. 4, pp. 1–12, 2018.
- [37] Y. Liu, B. Fan, S. Xiang, and C. Pan, "Relation-shape convolutional neural network for point cloud analysis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8895–8904.
- [38] S. Xie, S. Liu, Z. Chen, and Z. Tu, "Attentional shapecontextnet for point cloud recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4606–4615.
- [39] 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.
- [40] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, *et al.*, "Shapenet: An information-rich 3d model repository," *arXiv preprint arXiv:1512.03012*, 2015.
- [41] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black, "Smpl: A skinned multi-person linear model," *ACM transactions on graphics (TOG)*, vol. 34, no. 6, pp. 1–16, 2015.
- [42] A. Sharma, O. Grau, and M. Fritz, "Vconv-dae: Deep volumetric shape learning without object labels," in *European Conference on Computer Vision*. Springer, 2016, pp. 236–250.
- [43] J. Wu, C. Zhang, T. Xue, W. T. Freeman, and J. B. Tenenbaum, "Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling," in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, 2016, pp. 82–90.