

Introduction

This paper introduces APNet, an attention-based autoregressive network for task-oriented 3D point clouds sampling, which aims to sample a subset of points that are tailored specifically to a downstream task of interest. APNet employs a sequential autoregressive generation with a novel LSTM-based sequential model for sampling. Depending on the availability of labeled training data, APNet can be trained in supervised learning or self-supervised learning via knowledge distillation. We also present a joint training of APNet, yielding a single compact model that can generate arbitrary length of samples with prominent performances. Extensive experiments demonstrate the superior performance of APNet against state-of-the-arts in various downstream tasks, including 3D point cloud classification, reconstruction, and registration.

Method

Given original point cloud \mathbf{P} , the goal of APNet is to generate a point cloud $\mathbf{Q} = f_{\theta}(\mathbf{P})$ to maximize the predictive performance of task network \mathbf{T} . The parameters of APNet, θ , are optimized by minimizing a task loss and a sampling loss jointly as

$$\min_{\theta} \ell_{task}(T(\mathbf{Q}), y) + \lambda L_{sample}(\mathbf{Q}, \mathbf{P})$$

The sampling loss L_{sample} encourages the sampled points in \mathbf{Q} to be close to those of \mathbf{P} and also have a maximal coverage w.r.t. \mathbf{P} .

- Sampling loss

$$L_{sample}(\mathbf{Q}, \mathbf{P}) = L_a(\mathbf{Q}, \mathbf{P}) + \beta L_m(\mathbf{Q}, \mathbf{P}) + (\gamma + \delta |\mathbf{Q}|) L_a(\mathbf{P}, \mathbf{Q})$$

- Average nearest neighbor loss
- Maximal nearest neighbor loss

$$L_a(\mathbf{S}_1, \mathbf{S}_2) = \frac{1}{|\mathbf{S}_1|} \sum_{s_1 \in \mathbf{S}_1} \min_{s_2 \in \mathbf{S}_2} \|s_1 - s_2\|_2^2$$

$$L_m(\mathbf{S}_1, \mathbf{S}_2) = \max_{s_1 \in \mathbf{S}_1} \min_{s_2 \in \mathbf{S}_2} \|s_1 - s_2\|_2^2$$

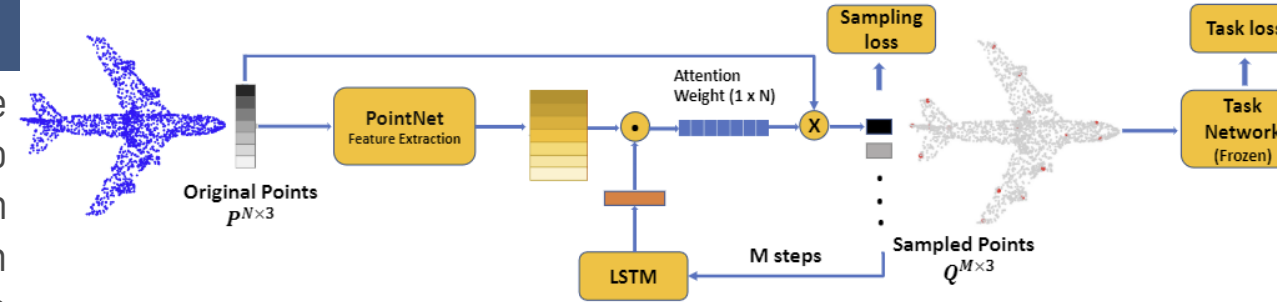


Fig. 1 Overview of APNet. APNet first extracts features with a simplified PointNet that preserves the geometric information of a point cloud. Then, an LSTM with attention mechanism is used to capture the relationship among points and select the most informative point sequentially. Finally, the sampled point cloud is fed to a task network for prediction. The whole pipeline is optimized by minimizing a task loss and a sampling loss jointly.

- Self-supervised Training with Knowledge Distillation
The task network \mathbf{T} is the teacher model, and APNet is the student model and use the soft predictions of \mathbf{T} as the targets to train APNet. In this scenario,

$$\ell_{task}(T(\mathbf{Q}), \tilde{y}), \text{ with } \tilde{y} = T(\mathbf{P})$$

- Joint Training
Given the autoregressive model of our method, APNet can generate arbitrary length of samples from a single mode. We can train one APNet with different sample sizes by

$$L_{joint} = \sum_{c \in C_s} (\ell_{task}(T(\mathbf{Q}_c), y) + \lambda L_{sample}(\mathbf{Q}_c, \mathbf{P}))$$

where C_s is a set of sample sizes of interest.

3. Inference Time

m	32	128	256	512
SampleNet-G	7.63	7.54	7.79	7.94
SampleNet-M*	44.33	135.23	261.47	515.30
APNet-G*	9.21	12.84	17.68	27.48
APNet-M	45.91	139.83	269.40	525.38

Experimental Results

1. Classification

m	RS	FPS	DaNet	MOPS-Net		SampleNet			APNet		APNet-KD	
				G	M	G	M	M*	G	M	G	M
8	8.26	23.29	-	-	-	78.36	73.31	28.7	81.42	74.12	80.22	73.81
16	25.11	54.19	-	84.7	51.2	80.60	79.68	55.5	83.89	82.25	83.82	82.02
32	55.19	77.32	85.1	86.1	77.6	80.32	82.97	74.4	88.15	86.97	88.76	84.95
64	78.26	87.22	86.8	87.1	81.0	79.36	84.01	79.0	88.38	87.58	88.66	87.54
128	85.95	88.76	86.8	87.2	85.0	85.52	87.17	79.7	89.22	89.38	87.83	88.01
256	88.80	89.30	87.2	87.4	86.7	87.43	89.58	83.4	89.54	89.86	88.02	88.21
512	89.66	89.87	-	88.3	88.3	88.01	90.18	88.2	89.78	90.18	88.69	88.56

Classification accuracies with different sample sizes m on ModelNet40.

2. Reconstruction

m	RS	FPS	SampleNet			APNet		APNet-KD	
			G	M	M*	G	M	G	M
8	21.85	12.79	5.29	5.48	-	4.27	4.59	4.69	4.98
16	13.47	7.25	2.78	2.89	-	2.51	2.62	2.57	2.67
32	8.16	3.84	1.68	1.71	2.32	1.54	1.59	1.47	1.52
64	4.54	2.23	1.32	1.27	1.33	1.07	1.11	1.12	1.14

The normalized reconstruction errors with different sample sizes m on the ShapeNet Core55 dataset.

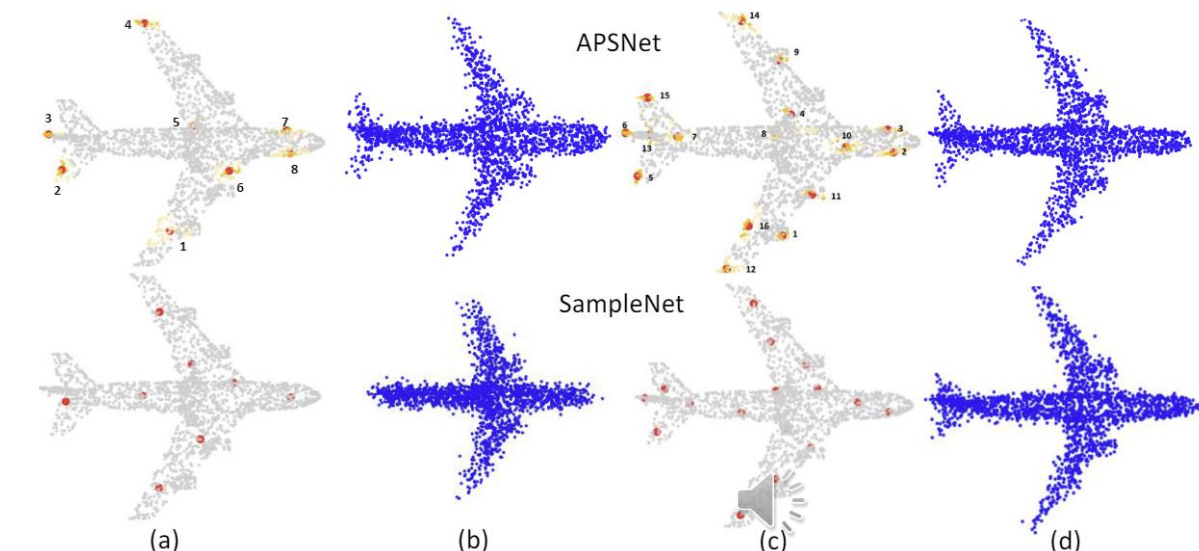


Fig. 2 Visualization of sampled points and reconstructed point clouds. APNet focuses more on the outline of the airplane without losing details, which are critical for the reconstruction