

Learning to Sample

Oren Dovrat*
Tel-Aviv University
Oren.Dovrat@gmail.com

Itai Lang*
Tel-Aviv University
Itai.Lang83@gmail.com

Shai Avidan
Tel-Aviv University
avidan@eng.tau.ac.il

Abstract

Processing large point clouds is a challenging task. Therefore, the data is often sampled to a size that can be processed more easily. The question is how to sample the data? A popular sampling technique is Farthest Point Sampling (FPS). However, FPS is agnostic to a downstream application (classification, retrieval, etc.). The underlying assumption seems to be that minimizing the farthest point distance, as done by FPS, is a good proxy to other objective functions.

We show that it is better to learn how to sample. To do that, we propose a deep network to simplify 3D point clouds. The network, termed S-NET, takes a point cloud and produces a smaller point cloud that is optimized for a particular task. The simplified point cloud is not guaranteed to be a subset of the original point cloud. Therefore, we match it to a subset of the original points in a post-processing step. We contrast our approach with FPS by experimenting on two standard data sets and show significantly better results for a variety of applications. Our code is publicly available¹

1. Introduction

Capturing 3D data is getting easier in recent years and there is a growing number of 3D shape repositories available online. This data can be represented in a variety of ways, including point clouds, multi-view images and voxel grids. A point cloud contains information only about the surface of a 3D object, while a grid based representation also holds data about free space, making the former much more efficient. However, processing a point cloud can be challenging, since it may contain a lot of data points. Reducing the number of points can be beneficial in many aspects, such as reduction of power consumption, computational cost and communication load, to name a few.

One naive approach to reduce the data load is to randomly sample a subset of points. Another approach, which

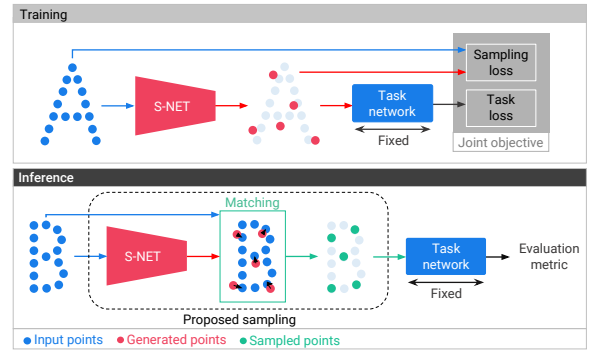


Figure 1. **An illustration of the proposed learned sampling approach.** In the *training* phase, S-NET generates points that are passed to a task network, which was pre-trained and is held fixed. The minimization objective contains the task’s loss and a sampling loss. The latter serves as a regularizer and encourages proximity between the input and generated points. At *inference* time, we match the points generated by S-NET with the input point cloud and get a subset of it. Only these points are then fed to the task network for performance evaluation.

is commonly used in the literature, is Farthest Point Sampling (FPS) [1, 2, 3]. This sampling method takes into account the structure of the point cloud and selects a group of points that are farthest apart from each other [4, 5]. These sampling methods, as well as other approaches in the literature [6, 7], operate according to a non-learned predetermined rule.

In the last few years, deep learning techniques have been applied with great success to point cloud data. Among various applications one can find point cloud classification [1, 2, 3, 8, 9, 10], part segmentation [1, 2, 3, 9, 10, 11], semantic segmentation [1, 3, 10, 12, 13, 14] and retrieval [15, 16]. Other techniques perform point cloud auto-encoding [17, 18, 19], generation [17, 20, 21], completion [17, 22, 23] and up-sampling [24, 25, 26]. Yet a learned point clouds sampling approach, subject to a subsequent task objective, has not been proposed before.

We propose a simplification network, termed S-NET, that is based on the architecture of PointNet [1]. S-NET

¹https://github.com/orendv/learning_to_sample

*Equal contribution

learns to generate a smaller (simplified) point cloud that is optimized for a downstream task, such as classification, retrieval or reconstruction.

The simplified point cloud must balance two conflicting constraints. On the one hand, we would like it to preserve similarity to the original shape. On the other hand, we wish to optimize it to a subsequent task. We solve this by training the network to generate a set of points that satisfy two objectives: a sampling loss and the task’s loss. The sampling loss drives the generated points close to the input point cloud. The task loss ensures that the points are optimal for the task.

An advantage of FPS is that it samples a subset of the original points. In contrast, the simplified point cloud produced by S-NET is not guaranteed to be a subset of the input point cloud. To address this issue, we perform a post-processing step at *inference* time, where we match the generated points with the input point cloud and obtain a subset of it, i.e., a set of sampled points (see Figure 1). Experiments show that better results for several tasks are achieved using our sampled points in comparison to FPS.

Our approach can be thought of as a feature selection mechanism [27, 28]. Each point is a feature of the underlying shape and we seek to select the ones that contribute the most to the task. It can also be interpreted as a form of visual attention [29, 30], focusing the subsequent task network on the significant points.

S-NET is trained to output a fixed sample size, which means that we need to train a different S-NET for every target size. To overcome this limitation, we introduce an extension of S-NET termed ProgressiveNet. ProgressiveNet orders points by importance to the task. This lets the sample size to be chosen at inference time, allowing for a dynamic level-of-detail management, according to the requirements and available resources.

The proposed sampling approach is applied to three different tasks: point cloud classification, retrieval and reconstruction. We compare our approach with common non-data driven methods: random sampling and FPS. For the first task we show better classification accuracy; in the second we show improved retrieval results; and in the last we get a lower reconstruction error. To summarize, our key contributions are:

- A task-specific data-driven sampling approach for point clouds;
- A Progressive sampling method that orders points according to their relevance for the task;
- Improved performance for point cloud classification, retrieval and reconstruction with sampled point clouds.

2. Related work

Point cloud simplification and sampling Several techniques for either point cloud simplification [31, 32] or sam-

pling [33, 34] have been proposed in the literature. Pauly *et al.* [31] presented and analyzed several simplification methods for point-sampled surfaces, including: clustering methods, iterative simplification and particle simulation. The simplified point set, resulting from these algorithms, was not restricted to be a subset of the original one. Farthest point sampling was adopted in the work of Moenning and Dodgson [32] as a means to simplify point clouds of geometric shapes, in a uniform as well as feature-sensitive manner.

Katz and Tal [33] suggested a view dependent algorithm to reduce the number of points. They used hidden-point removal and target-point occlusion operators in order to improve a human comprehension of the sampled point set. Recently, Chen *et al.* [34] employed graph-based filters to extract per point features. Points that preserve specific information are likely to be selected by their sampling strategy. The desired information is assumed to be beneficial to a subsequent application.

The above sampling approaches aim to optimize a variety of sampling objectives. However, they do not consider *directly* the objective of the task to be followed.

Progressive simplification In a seminal paper, Hoppe [35] proposed a technique for progressive mesh simplification. In each step of his method, one edge is collapsed such that minimal geometric distortion is introduced.

A recent work by Hanocka *et al.* [36] suggested a neural network that performs task-driven mesh simplification. Their network relies on the edges between the mesh vertices. This information is not available for point clouds.

Several researchers studied the topic of point set compression [37, 38, 39]. An octree data structure was used for progressive encoding of the point cloud. The objective of the compression process was low distortion error.

Deep learning on point sets The pioneering work of Qi *et al.* [1] presented PointNet, the first neural network that operates directly on unordered point cloud data. They constructed their network from per-point multi-layer perceptrons, a symmetric pooling operation and several fully connected layers. PointNet was employed for classification and segmentation tasks and showed impressive results. For assessing the applicability of PointNet for reduced number of input points, they used random sampling and FPS. In our work, we suggest a data-driven sampling approach, that improves the classification performance with sampled sets in comparison to these sampling methods.

Later on, Qi *et al.* extended their network architecture for hierarchical feature learning [2]. In the training phase, centroid points for local feature aggregation were selected by FPS. Similar to their previous work [1], FPS was used for evaluating the ability of their network to operate on fewer input points.

Li *et al.* [11] suggested to learn the centroid points for feature aggregation by a self-organizing map (SOM). They used feature propagation between points and SOM nodes and showed improved results for point cloud classification and part segmentation. The SOM was optimized separately, as a pre-processing step.

Building on the work of Qi *et al.* [1], Achlioptas *et al.* [17] developed autoencoders and generative adversarial networks for point clouds. Instead of shape class or per-point label, the output of their network was a set of 3D points. In this work we apply our sampling approach for the task of point cloud reconstruction with the autoencoder proposed by Achlioptas *et al.*

Several researchers tackled the problem of point cloud consolidation. Yu *et al.* [24] extracted point clouds from geodesic patches of mesh models. They randomly sampled the point sets and trained a network to reconstruct the original patch points. Their follow-up work [25] incorporated edge information to improve the reconstruction accuracy. Zhang *et al.* [26] studied the influence of sampling strategy on point cloud up-sampling. They used Monte-Carlo random sampling and curvature based sampling. Their network was trained to produce a fixed size point cloud from its sample. In contrast to these works, our study focuses on the down-sampling strategy.

3. Method

Problem statement Given a point set $P = \{p_i \in \mathbb{R}^3, i = 1, \dots, n\}$, a sample size $k \leq n$ and a task network T , find a subset S^* of k points that minimizes the task network’s objective function f :

$$S^* = \underset{S}{\operatorname{argmin}} f(T(S)), \quad S \subset P, \quad |S| = k \leq n. \quad (1)$$

This problem poses a challenge, as sampling might seem akin to pooling, yet in pooling the pooled value is propagated forward, so the gradient with respect to it can be calculated. Discrete sampling, however, is like “arg-pooling”, where the propagated value cannot be updated incrementally. As a result, a sampling operation cannot be trained directly. Therefore, we propose a two-step process: first, we employ a neural network, i.e., S-NET, to generate a set of points. Second, we match the *generated* points with the input point cloud to obtain a subset of its points, i.e., the *sampled* points. Figure 1 illustrates the process.

The input to S-NET is a set of n 3D coordinates, namely points, representing a 3D shape. The output of S-NET is k generated points. S-NET is followed by a task network. The task network is pre-trained on an input of n points, to perform a given task on the point cloud (i.e., classification, retrieval or reconstruction). It is kept fixed during training and testing of S-NET. This ensures that sampling is being optimized to the task, rather than the task being optimized to an arbitrary sampling.

At the training phase, the generated points are fed to the task network. The points are optimized to the task at hand by minimizing the task loss. We use an additional sampling regularization loss term, that encourages each of the generated points to be close to one of the input points and forces the generated points to spread over the input cloud.

At inference, the generated points are matched with the input point cloud in order to obtain a subset of it. These are the sampled points, the final output of our process. These points are passed through the task network and its performance is evaluated.

We present two sampling versions: S-NET and ProgressiveNet. In the first version (Figure 1), we train a different sampling network per sample size. In the second one (Figure 2), we train one network that can be used to produce any sample size smaller than the input size.

3.1. S-NET

The architecture of S-NET follows that of Qi *et al.* [1]. The input points undergo a set of 1×1 convolution layers, resulting in a per point feature vector. Then, a symmetric feature-wise max pooling operation is used to obtain a global feature vector. Finally, we use several fully-connected layers. The output of the last layer is the set of generated points.

Let us denote the generated point set as G and the input point set as P . We construct a sampling regularization loss, composed out of three terms:

$$L_f(G, P) = \frac{1}{|G|} \sum_{g \in G} \min_{p \in P} \|g - p\|_2^2 \quad (2)$$

$$L_m(G, P) = \max_{g \in G} \min_{p \in P} \|g - p\|_2^2 \quad (3)$$

$$L_b(G, P) = \frac{1}{|P|} \sum_{p \in P} \min_{g \in G} \|p - g\|_2^2. \quad (4)$$

L_f and L_m keeps the points in G close to those in P , in the average and worst case, respectively. This is designed to encourage tight matches in the following matching process. We found that mixing average and maximum operations speeds up convergence. L_b ensures that the generated points are well spread over the input points, decreasing the number of collisions in the matching process. The sampling regularization loss is a weighted sum of these three terms:

$$L_s(G, P) = L_f(G, P) + \beta L_m(G, P) + (\gamma + \delta |G|) L_b(G, P). \quad (5)$$

Note that this is a generalization of the Chamfer distance [40], achieved when $\beta = 0$, $\gamma = 1$ and $\delta = 0$.

In addition, we denote L_{task} as the task network loss. The total S-NET loss is:

$$L^{S-NET}(G, P) = L_{task}(G) + \alpha L_s(G, P) \quad (6)$$

where α controls the regularization trade-off. The output of S-NET is a $k \times 3$ matrix, where k is the sample size, and we train separately for each k .

3.2. Matching

The generated points G are not guaranteed to be a subset of the input points P . In order to get a subset of the input points, we match the generated points to the input point cloud.

A widely used approach for matching two point sets is the Earth Mover’s Distance (EMD) [17, 23, 26, 40]. EMD finds a bijection between the sets that minimizes the average distance of corresponding points, while the point sets are required to have the same size. In our case, however, G and P are of different size.

We examine two matching methods. The first adapts EMD to uneven point sets. The second is based on nearest neighbour (NN) matching. Here we describe the latter, which yielded better results. The reader is referred to the supplementary material for details about the other matching method.

In NN-based matching, each point $x \in G$ is replaced with its closest euclidean corresponding point $y^* \in P$:

$$y^* = \operatorname{argmin}_{y \in P} \|x - y\|_2. \quad (7)$$

Since several points in G might be closest to the same point in P , the number of unique sampled points might be smaller than the requested sample size. Therefore, we remove duplicate points and get an initial sampled set. Then, we complete this set, up to the size of G , by running farthest point sampling (FPS) [2], where in each step we add a point from P that is farthest from the current sampled set.

The matching process is only applied at inference time, as the final step of inference. During training, the generated points are processed by the task network as-is, since the matching is not differentiable and cannot propagate the task loss back to S-NET.

3.3. ProgressiveNet: sampling as ordering

S-NET is trained to sample the points to a single, predefined, sample size. If more than one sample size is required, more than one S-NET needs to be trained. But what if we want to train one network that can produce any sample size, i.e., sample the input at any sampling ratio? To this end we present ProgressiveNet. ProgressiveNet is trained to take a point cloud of a given size and return a point cloud of the same size, consisting of the same points. But, while the points of the input are arbitrarily ordered, the points of the output are ordered by their relevance to the task. This allows sampling to any sample size: to get a sample of size k , we simply take the first k points of the output point cloud of ProgressiveNet and discard the rest. The architecture of

ProgressiveNet is the same as S-NET, with the last fully connected layer size equal to the input point cloud size (has $3n$ elements).

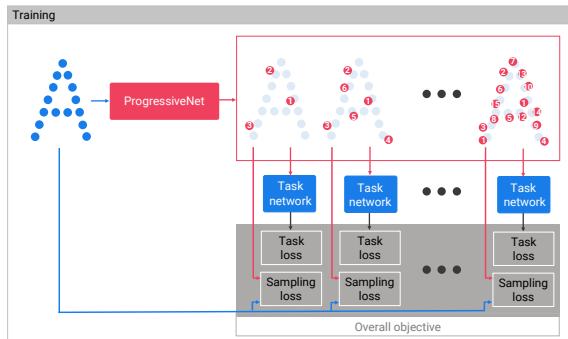


Figure 2. **ProgressiveNet training.** The generated points are divided into groups of increasing size, such that each group is a subset of the following larger group. Each group has corresponding task and sampling losses. The overall objective is the sum of the per-group losses. The task network was pre-trained and is kept fixed.

To train ProgressiveNet (Figure 2) we define a set of sizes $C_s = \{2^1, 2^2, \dots, 2^{\log_2(n)}\}$. For each size $c \in C_s$ we compute a task loss term and a sampling regularization loss term, such that the total ProgressiveNet’s loss becomes:

$$L^{ProgressiveNet}(G, P) = \sum_{c \in C_s} L^{S-NET}(G_c, P) \quad (8)$$

where L^{S-NET} is the loss term as defined in equation 6 and G_c are the first c points from the points generated by ProgressiveNet (the first $3c$ elements of the output layer). This loss function defines a nested structure of point subsets. For example, the first subset, consisting of two points, is used in all terms $L^{S-NET}(G_c, P)$, for $c \geq 2$. Because this subset is used in all terms, it is contained in all larger subsets. Similarly, the first 4 points define the second subset, that includes the first subset and is part of all larger subsets.

Under this training process, the first k generated points (for any k) are optimized to be suitable both for the task at hand as a stand-alone set and for integration with their preceding points to create a larger set that will improve results on the given task. This makes sure that a point that is more important for the task will appear earlier in the generated point set, while extra points will give diminishing marginal utility, resulting in a task-oriented progressive decomposition of the point cloud [41].

At inference, the generated points (ProgressiveNet’s output layer) are matched with the input point cloud using the same matching process we used in Section 3.2 for S-NET. We transfer the order of the generated points to their matched points. To obtain a specific sample size k , we take the first k sampled points.

4. Results

We applied S-NET to three tasks: point set classification [1], retrieval and reconstruction [17]. For classification and retrieval we adopt the ModelNet40 [42] point cloud data provided by Qi *et al.* [1] and PointNet [1] as the task network. For reconstruction we employ ShapeNet Core55 repository [43] and the point set autoencoder of Achliopas *et al.* [17].

Random sampling and FPS are employed as alternative non-data driven sampling approaches for comparison with our suggested approach. In order to make a fair comparison, we only use S-NET points after the matching process, so both our sampled points and the alternative approaches’ points are subsets of the input point cloud. Further experimental details can be found in the supplementary material.

4.1. Classification

We report instance classification results on the ModelNet40 data set, adopting the official train-test split. We employed the full version of PointNet as the task network for S-NET and the vanilla version of PointNet for ProgressiveNet (to save training time), except where otherwise noted.

S-NET Figure 3 shows the classification accuracy of PointNet, when trained on the complete data (1024 points per point cloud) and tested on samples of different size. We compare several different sampling methods: random sampling with uniform distribution over the input points, FPS and S-NET. We trained 10 different S-NETs, for sample sizes of $k \in \{2, 4, 8, \dots, 1024\}$ points. The sampling ratio is defined as $n/k = 1024/k$. We observe that the classification accuracy using S-NET’s sampled points is equal or better than that of using FPS’s points for any sampling ratio, with a margin of up to 34.2% (for sampling ratio 32).

Table 1 shows that the accuracy of PointNet is also higher when *training* it on the points sampled by S-NET. Each of the numbers in this table represents a different PointNet, each trained and tested on point clouds of a specific size that was sampled with a different sampling method. We see, for example, that PointNet that was both trained and tested on point clouds of just 16 points, achieved 85.6% accuracy, when using S-NET, compared to only 76.7% with FPS. This shows that S-NET is not overfitted to the specific instance of the classifier it was trained with. Instead, it samples the points in a way that makes the sampled set easy to classify, creating a strong distinction between different classes in the data.

ProgressiveNet In Figure 4 we compare the accuracy of PointNet vanilla on points sampled by S-NET (trained with PointNet vanilla, in this case) to those sampled by ProgressiveNet. The evaluation was done for all sample sizes in

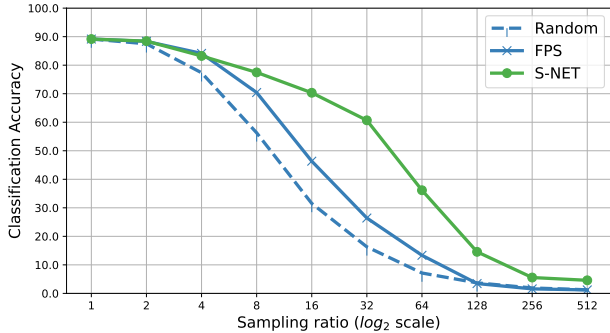


Figure 3. **S-NET for classification.** PointNet was trained on complete point clouds (1024 points) and evaluated on sampled point clouds of the test set using different sampling methods: random, FPS, and S-NET. The accuracy using S-NET is evidently higher.

#Sampled points	Random	FPS	S-NET
1024	89.2	89.2	89.2
512	88.2	88.3	87.8
256	86.6	88.1	88.3
128	86.2	87.9	88.6
64	81.5	86.1	87.7
32	77.0	82.2	87.3
16	65.8	76.7	85.6
8	45.8	61.6	83.6
4	26.9	35.2	73.4
2	16.6	18.3	53.0

Table 1. **Training PointNet classifier on sampled point clouds.** We sample the train and test data with different sampling methods: random sampling, FPS, and S-NET. Afterwards, PointNet is trained and evaluated on the sampled data. Applying S-NET allows good classification even with minimal data to train on. We conclude that S-NET transforms the data into a more separable representation.

the range [2, 1024]. S-NET results are from 10 different S-NETs, each trained for a specific sample size, for sizes of $k \in \{2, 4, 8, \dots, 1024\}$ points. For the sample sizes in between those values, we took the S-NET that was trained for a lower sample size and then completed with FPS to the required size. For example, to get a sample of size 48, we took the points sampled by S-NET that was trained to sample 32 points, and then made 16 steps of FPS. The Progressive results are from one ProgressiveNet with output size 1024, that was trained with classification and sampling loss terms for sizes $C_s = \{2, 4, 8, \dots, 1024\}$.

We observe that S-NET has better performance for the sample sizes it was trained for, while ProgressiveNet performs better for any sample size in between. This shows the advantage of ProgressiveNet, which orders the points by priority, so that the accuracy is approximately monotonically increasing in the sample size.

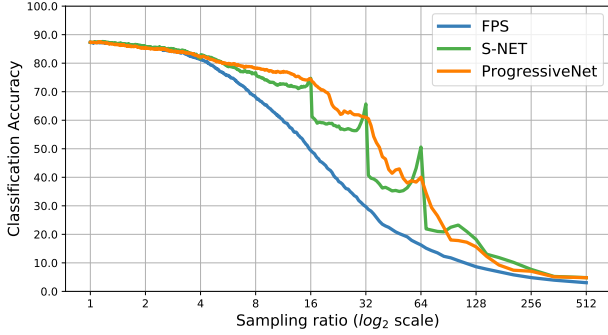


Figure 4. **ProgressiveNet vs. S-NET.** PointNet vanilla was trained on complete point clouds (1024 points) and evaluated on every sample size k in the range $[2, 1024]$. The sampling ratio is defined as $1024/k$. We compared three different sampling methods: FPS, S-NET and ProgressiveNet. S-NET is better for the sample sizes it was trained for, while ProgressiveNet performs better for sample sizes in between.

Scalability by S-NET S-NET assumes that the task network is already trained. This will cause a problem in case of very large point clouds. We show that it is possible to train the task network on FPS-sampled point clouds and use that to train S-NET. The resulting S-NET can then be used to re-train the task network to a higher accuracy. Please see Figure 5 for an illustration of the proposed training and testing procedure.

The following small scale simulation demonstrate this: We used FPS to sample ModelNet40 to $k = 32$ points per point cloud, and trained PointNet on those sampled points. This is our baseline. The accuracy of the baseline on the test set is 82.2%. When we fed this trained PointNet with point clouds of size 1024, the accuracy was just 46.6%. We then trained S-NET, using the 32-sized point clouds training set, and the baseline PointNet as the task network. Then we sampled ModelNet40 again to size $k = 32$, this time using S-NET. Finally, we trained PointNet on the points sampled by S-NET. The accuracy of the re-trained PointNet on the test set improved to 86.0%. Employing S-NET allows us to improve PointNet’s accuracy without training on larger point clouds.

Time and space considerations The time complexity of PointNet-like networks is dominated by their per-point convolution layers, and thus is strongly dependent on the size of the input point cloud. The space complexity of S-NET is a linear function of its output size k . S-NET offers a trade-off between space (number of parameters) and inference time (number of floating point operations). For example, cascading S-NET that samples a point cloud of 1024 to 16 points with a following PointNet reduces inference time by over 90% compared to running PointNet on the complete point cloud, with only 5% increase in space. See full details in the supplementary material.

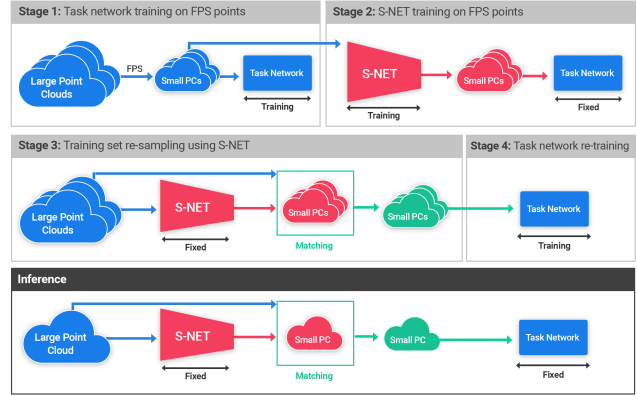


Figure 5. **Scalability by S-NET.** Illustration of the proposed training and inference procedure on large point clouds. In stage 1 we use FPS to sample the point clouds to a trainable size and use the sampled point clouds to train the task network. In stage 2 we train S-NET on the sampled point clouds, employing the fixed task network. In stage 3 we apply S-NET to re-sample the large point clouds and in stage 4 we train the task network again, this time on S-NET’s points. At inference time, we use S-NET to sample a large point cloud to the size the task network was trained for.

Approximate sampling Up to this point we applied the matching post-processing step to compare ourselves directly with FPS. However, there might be settings where the k output points do not have to be a subset of the original point cloud. In such cases, we can use S-NET’s generated points directly, forgoing the matching step. One can use either the generated or the sampled points. A third alternative is to interpolate between the two, i.e., using points that are up to ϵ away from the original input points. This is done by interpolating each generated point with its matched sampled point, to get a third point on the line between them, no more than ϵ away from the sampled point. Figure 6 shows that PointNet’s accuracy is higher when feeding it the generated points. For point clouds normalized to the unit sphere, we find that choosing $\epsilon = 0.05$ results in classification accuracy that is about mid-way between the accuracy on the sampled and generated points. Note that ϵ is set at inference time.

Critical set sampling The critical set, as defined by Qi *et al.* [1], is the set of points that contributed to the max pooled features. A plausible alternative to our method might be to sample the critical points that contribute the most features. We tried this approach and found it viable only at small sampling ratios. See full details in the supplementary material.

4.2. Retrieval

We now show that employing S-NET instead of FPS leads to better retrieval results. Specifically, we took the S-NET that was trained with PointNet for classification as the task network, without re-training it. We sampled the points

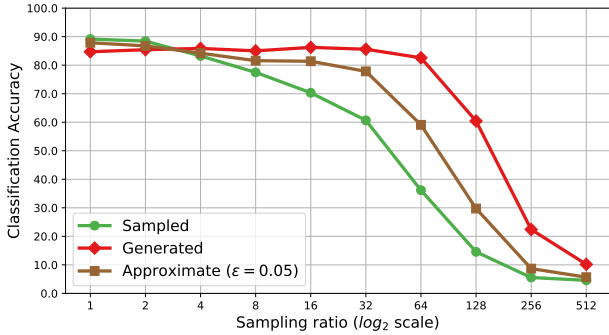


Figure 6. **Approximate sampling.** PointNet was trained on complete point clouds (1024 points) and evaluated on point clouds of different sizes. We use S-NET’s generated and sampled points, as well as a third set of points, where each point is an interpolation between the generated and sampled point, bounded to be no more than $\epsilon = 0.05$ away from an original input point. Approximate sampling enables higher accuracy, when deviating from original points is not a concern.

Sampling ratio	FPS mAP	S-NET mAP
1	71.3	71.3
2	70.1	69.8
4	65.7	64.8
8	58.3	60.4
16	49.4	59.0
32	37.7	59.0
64	27.4	54.5

Table 2. **Point cloud retrieval.** We took the same S-NET that was trained with PointNet classifier as the task network and applied it as the sampling method for retrieval. The shape descriptor was PointNet’s activations of the layer before the last, with L_2 as the distance metric. We measured the macro mean Average Precision (mAP) for different sampling ratios and methods. S-NET performs better than FPS for large sampling ratios and is almost insensitive to the sampling ratio.

that are fed to PointNet and used its penultimate layer as a shape descriptor. Retrieval was done based on L_2 distance on this shape descriptor. We repeated the experiment with every shape in the ModelNet40 test set serving as a query shape. The experiment was repeated with different sample sizes for both S-NET and FPS.

Table 2 summarizes the mean Average Precision (mAP) for different sampling ratios and sampling methods. We see that S-NET performs much better for sampling ratios larger than 4 and is not very sensitive to the sampling ratio. Figure 7 presents the precision-recall curve for the original data and for the sampled data using FPS and S-NET, sampling to $k=32$ points per shape. We observe significantly better precision when applying S-NET across all recall values.

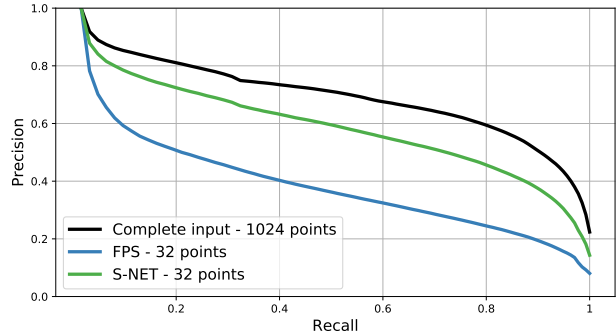


Figure 7. **Precision-Recall curve for point cloud retrieval.** We compare FPS and S-NET when sampling 32 points, as well as using the complete 1024 points data. Significantly better precision is achieved when using S-NET compared to FPS, across all recall values.

4.3. Reconstruction

We next learn to sample with an autoencoder as the task network. We used the ShapeNet Core55 point cloud data provided by Achlioptas *et al.* [17], as well as their autoencoder and trained it on four shape classes: table, car, chair and airplane. These classes have the most available models. Each shape class is split into 85%-5%-10% for train-validation-test sets. The autoencoder was trained to receive and reconstruct point clouds of 2048 points.

The reconstruction error of the autoencoder is measured by the Chamfer distance [17]. In order to compare different sampling methods, we use Normalized Reconstruction Error (NRE). That is, we reconstruct the complete point set from a subset of points and from the complete point set and take the reconstruction error ratio between the two. We trained several S-NETs with the following sample sizes: $k \in \{16, 32, \dots, 2048\}$, and a single ProgressiveNet with loss terms for the same sizes. As an alternative sampling approach we used FPS.

Normalized reconstruction error Figure 8 presents the NRE as a function of the sampling ratio, where the sampling ratio is defined as $n/k = 2048/k$. We compare FPS with S-NET and ProgressiveNet. For small sampling ratios, the NRE for our sampled points is similar to that of FPS. However, as the sampling ratio increases, our sampling methods outperform FPS. For example, at sampling ratio of 32, the NRE of FPS is a little over 2, while the NRE of S-NET and ProgressiveNet is about 1.5 and 1.75, respectively. S-NET achieves lower NRE than ProgressiveNet, since the former was optimized separately per sampling ratio, resulting in improved reconstruction performance. We learn to sample points from unseen shapes that enable lower reconstruction error.

Sample and reconstruction visualization Figure 9 compares the reconstruction result from the entire point cloud,

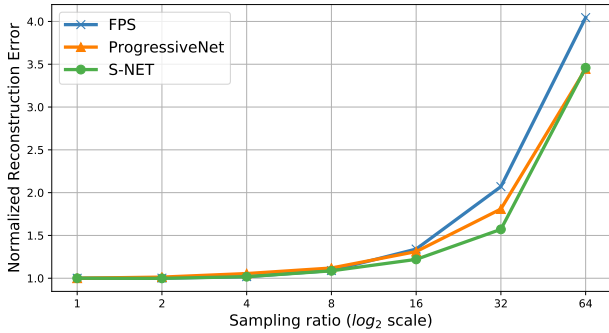


Figure 8. **Normalized Reconstruction Error (NRE)**. We trained an autoencoder on complete point clouds of 2048 points, and evaluated the reconstruction error from sampled point clouds with FPS, S-NET and ProgressiveNet on the test split. Up to a sampling ratio of 8, the error for S-NET and ProgressiveNet is on par with FPS. However, at higher sampling ratios S-NET and ProgressiveNet achieves lower error.

to that from a sample size of 64 points, where the samples are produced by either S-NET or FPS. The reconstruction quality when employing S-NET is higher than that of using FPS and approaches that of using the entire point cloud. Interestingly, the points sampled by S-NET are *non-uniformly* distributed, as opposed to the more uniform distribution of the points sampled by FPS.

Adversarial simplification In this proof of concept we show how to trick the autoencoder. We simplify a point cloud to be visually similar to one class but reconstructed by the autoencoder to a shape from a *different* class. We train S-NET with a single pair of input and target shapes, where the input is a shape from one class and the target is a shape from another class. The sampling loss was between the input and the points generated by S-NET. The reconstruction loss was between the target shape and the reconstructed one. Figure 10 shows the result of turning an airplane into a car.

5. Conclusions

We presented a method that learns how to sample a point cloud that is optimized for a downstream task. The method consists of a simplifying network, S-NET, followed by a post-processing matching step. We also suggested a network, termed ProgressiveNet, for ordering a point cloud according to the contribution of each point to the task. The resulting methods outperform FPS in sampling points for several tasks: classification, retrieval and reconstruction of point clouds.

The proposed method is general and can produce a small point cloud that consists of points that are not necessarily part of the original input shape. The output point cloud minimizes a geometric error with respect to the input point

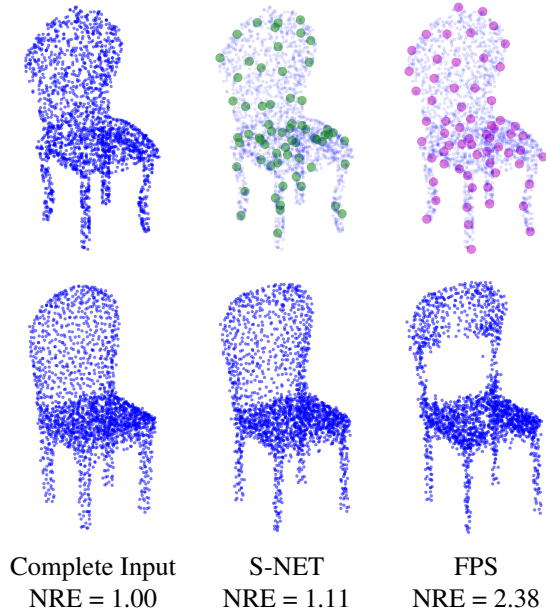


Figure 9. **Point cloud reconstruction**. NRE stands for Normalized Reconstruction Error. Top row: complete input point cloud of 2048 points, input with 64 S-NET sampled points (in Green), input with 64 FPS points (in Magenta). The sampled and FPS points are enlarged for visualization purpose. Bottom row: reconstructed point cloud from the input and from the corresponding sample. The reconstructed point cloud from S-NET’s sampled points is visually more similar to the input and has lower reconstruction error.

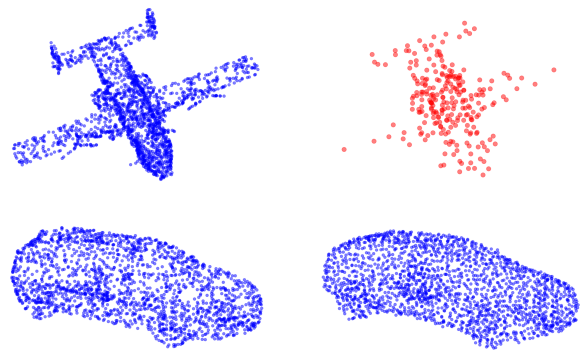


Figure 10. **Adversarial simplification**. Top row: input shape (in Blue) and 256 generated points (in Red). Bottom row: target shape and reconstruction from the generated points. While the simplified point cloud resembles the input airplane shape, it is reconstructed to a *completely different* shape - a car!

cloud while optimizing the objective function of a downstream task. We have shown that learning to sample can improve results and be used in conjunction with various applications.

Acknowledgments: Parts of this research were supported by ISF grant 1917/15.

References

- [1] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 652–660, 2017. [1](#), [2](#), [3](#), [5](#), [6](#), [11](#), [13](#)
- [2] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space," *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, 2017. [1](#), [2](#), [4](#)
- [3] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, "PointCNN: Convolution On X-Transformed Points," *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, 2018. [1](#)
- [4] Y. Eldar, M. Lindenbaum, M. Porat, and Y. Y. Zeevi, "The Farthest Point Strategy for Progressive Image Sampling," *IEEE Transactions on Image Processing*, vol. 6, pp. 1305–1315, 1997. [1](#)
- [5] C. Moenning and N. A. Dodgson, "Fast Marching farthest point sampling," *Eurographics Poster Presentation*, 2003. [1](#)
- [6] M. Alexa, J. Behr, D. Cohen-Or, S. Fleishman, D. Levin, and C. T. Silva, "Point Set Surfaces," *Proceedings of IEEE Visualization Conference*, pp. 21–28, 2001. [1](#)
- [7] L. Lars, "Point Cloud Representation," *Technical Report, Faculty of Computer Science, University of Karlsruhe*, 2001. [1](#)
- [8] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Pczos, R. Salakhutdinov, and A. J. Smola, "Deep Sets," *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, 2017. [1](#)
- [9] Y. Shen, C. Feng, Y. Yang, and D. Tian, "Mining Point Cloud Local Structures by Kernel Correlation and Graph Pooling," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [1](#)
- [10] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic Graph CNN for Learning on Point Clouds," *arXiv preprint arXiv:1801.07829*, 2018. [1](#)
- [11] J. Li, B. M. Chen, and G. H. Lee, "SO-Net: Self-Organizing Network for Point Cloud Analysis," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [1](#), [3](#)
- [12] L. Tchapmi, C. Choy, I. Armeni, J. Gwak, and S. Savarese, "SEGCloud: Semantic Segmentation of 3D Point Clouds," *Proceedings of the International Conference on 3D Vision (3DV)*, 2017. [1](#)
- [13] H. Su, V. Jampani, D. Sun, S. Maji, E. Kalogerakis, M.-H. Yang, and J. Kautz, "SPLATNet: Sparse Lattice Networks for Point Cloud Processing," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2530–2539, 2018. [1](#)
- [14] B.-S. Hua, M.-K. Tran, and S.-K. Yeung, "Pointwise Convolutional Neural Networks," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [1](#)
- [15] M. A. Uy and G. H. Lee, "PointNetVLAD: Deep Point Cloud Based Retrieval for Large-Scale Place Recognition," *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [1](#)
- [16] Z. Kuang, J. Yu, J. Fan, and M. Tan, "Deep Point Convolutional Approach for 3D Model Retrieval," *Proceedings of the IEEE International Conference on Multimedia and Expo (ICME)*, 2018. [1](#)
- [17] P. Achlioptas, O. Diamanti, I. Mitliagkas, and L. J. Guibas, "Learning Representations and Generative Models For 3D Point Clouds," *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018. [1](#), [3](#), [4](#), [5](#), [7](#), [16](#)
- [18] Y. Yang, C. Feng, Y. Shen, and D. Tian, "FoldingNet: Point Cloud Auto-encoder via Deep Grid Deformation," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [1](#)
- [19] T. Groueix, M. Fisher, V. G. Kim, B. Russell, and M. Aubry, "AtlasNet: A Papier-Mâché Approach to Learning 3D Surface Generation," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [1](#)
- [20] Y. Sun, Y. Wang, Z. Liu, J. E. Siegel, and S. E. Sarma, "PointGrow: Autoregressively Learned Point Cloud Generation with Self-Attention," *arXiv preprint arXiv:1810.05591*, 2018. [1](#)
- [21] C.-L. Li, M. Zaheer, Y. Zhang, B. Poczos, and R. Salakhutdinov, "Point Cloud GAN," *arXiv preprint arXiv:1810.05795*, 2018. [1](#)
- [22] S. Agrawal and S. Gurumurthy, "High Fidelity Semantic Shape Completion for Point Clouds using Latent Optimization," *Proceedings of the IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2018. [1](#)
- [23] W. Yuan, T. Khot, D. Held, C. Mertz, and M. Hebert, "PCN: Point Completion Network," *Proceedings of the International Conference on 3D Vision (3DV)*, 2018. [1](#), [4](#)
- [24] L. Yu, X. Li, C.-W. Fu, D. Cohen-Or, and P. A. Heng, "PU-Net: Point Cloud Upsampling Network," *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [1](#), [3](#)
- [25] L. Yu, X. Li, C.-W. Fu, D. Cohen Or, and P. A. Heng, "EC-Net: an Edge-aware Point set Consolidation Network," *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. [1](#), [3](#)
- [26] W. Zhang, H. Jiang, Z. Yang, S. Yamakawa, K. Shimada, and L. B. Kara, "Data-driven Upsampling of Point Clouds," *arXiv preprint arXiv:1807.02740*, 2018. [1](#), [3](#), [4](#)
- [27] I. Guyon and A. Elisseeff, "An introduction to variable and feature selection," *The Journal of Machine Learning Research*, pp. 1157–1182, 2003. [2](#)
- [28] R. Kohavi and G. H. John, "Wrappers for feature subset selection," *Artificial Intelligence*, pp. 273–324, 1997. [2](#)
- [29] V. Mnih, N. Heess, A. Graves, and K. Kavukcuoglu, "Recurrent Models of Visual Attention," *Proceedings of Advances in Neural Information Processing Systems (NIPS)*, 2014. [2](#)

- [30] B. Jimmy, V. Mnih, and K. Kavukcuoglu, “Multiple Object Recognition with Visual Attention,” *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015. 2
- [31] M. Pauly, M. Gross, and L. P. Kobbelt, “Efficient Simplification of Point-Sampled Surfaces,” *Proceedings of IEEE Visualization Conference*, 2002. 2
- [32] C. Moenning and N. A. Dodgson, “A new point cloud simplification algorithm,” *Proceedings of the IASTED International Conference on Visualization, Imaging and Image Processing (VIIP)*, 2003. 2
- [33] S. Katz and A. Tal, “Improving the Visual Comprehension of Point Sets,” *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 121–128, 2013. 2
- [34] S. Chen, D. Tian, C. Feng, A. Vetro, and J. Kovaevi, “Fast Resampling of Three-Dimensional Point Clouds via Graphs,” *IEEE Transactions on Signal Processing*, vol. 66, pp. 666–681, 2018. 2
- [35] H. Hoppe, “Progressive Meshes,” *Proceedings of the ACM Special Interest Group on Computer Graphics (SIGGRAPH)*, pp. 99–108, 1996. 2
- [36] R. Hanocka, A. Hertz, N. Fish, R. Giryes, S. Fleishman, and D. Cohen-Or, “MeshCNN: A Network with an Edge,” *arXiv preprint arXiv:1809.05910*, 2018. 2
- [37] J. Peng and C.-C. J. Kuo, “Octree-Based Progressive Geometry Encoder,” *Proceedings of SPIE*, pp. 301–311, 2003. 2
- [38] Y. Huang, J. Peng, C.-C. J. Kuo, and M. Gopi, “Octree-Based Progressive Geometry Coding of Point Clouds,” *Proceedings of the Eurographics Symposium on Point-Based Graphics*, pp. 103–110, 2006. 2
- [39] R. Schnabel and R. Klein, “Octree-based Point-Cloud Compression,” *Proceedings of the Eurographics Symposium on Point-Based Graphic*, 2006. 2
- [40] H. Fan, H. Su, and L. J. Guibas, “A point set generation network for 3d object reconstruction from a single image,” *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 3, 4, 11
- [41] J. M. Singh and P. Narayanan, “Progressive Decomposition of Point Clouds Without Local Planes,” *Computer Vision, Graphics and Image Processing*, pp. 364–375, 2006. 4
- [42] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, “3D ShapeNets: A Deep Representation for Volumetric Shapes,” *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1912–1920, 2015. 5
- [43] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu, “ShapeNet: An Information-Rich 3D Model Repository,” *arXiv preprint arXiv:1512.03012*, 2015. 5
- [44] V. Nair and G. E. Hinton, “Rectified Linear Units Improve Restricted Boltzmann Machines,” *Proceedings of the International Conference on Machine Learning (ICML)*, 2010. 11
- [45] S. Ioffe and C. Szegedy, “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift,” *Proceedings of the International Conference on Machine Learning (ICML)*, 2015. 11
- [46] D. P. Kingma and J. L. Ba, “Adam: A Method for Stochastic Optimization,” *Proceedings of the International Conference on Learning Representations (ICLR)*, 2015. 11
- [47] D. P. Bertsekas, “A distributed asynchronous relaxation algorithm for the assignment problem,” *Proceedings of the 24th IEEE Conference on Decision and Control*, pp. 1703–1704, 1985. 11

Supplementary

In the following sections we provide additional details and analysis regarding our results. In Section A we provide experimental details and additional results for our method. Section B displays visual examples of retrieved shapes and progressive sampling. Section C offers a new application for the progressive concept: A progressive autoencoder.

A. Additional results

In this section we extend the results section of the paper to include: details of our experimental parameters, analysis of the different matching methods, evaluation of regularization parameters, elaboration of the space and time considerations, analysis of critical set sampling, and comparison of the class accuracy.

A.1. Experimental details

Classification Architecture The architecture of S-NET for the classification experiment is inspired by the vanilla version of PointNet by Qi *et al.* [1]. We use per point 1×1 convolution layers with output sizes of [64, 64, 64, 128, 128]. Then, a feature-wise max-pooling layer is used to obtain a global feature vector. This feature vector is passed through four fully connected layers of size [256, 256, 256, $k \times 3$], where k is the sample size.. All convolution and fully connected layers are followed by ReLU non-linearity [44] and batch-normalization layer [45], except for the output layer. ProgressiveNet takes the architecture of S-NET with $k = 1024$.

Classification Optimization We used Adam optimizer [46] with a learning rate of 0.01, decay rate of 0.7 every 60000 steps and batch size of 32 point clouds. The regularization weights (for equations 5 and 6) were set to $\alpha = 30$, $\beta = 1$, $\gamma = 1$, $\delta = 0$ for S-NET and $\alpha = 30$, $\beta = 1$, $\gamma = 0.5$, $\delta = 1/30$ for ProgressiveNet. We trained the networks for 500 epochs with a GTX 1080 Ti GPU. When using PointNet as the task network, S-NET training takes between 2 to 7 hours, depending on the sample size. ProgressiveNet training takes 6 hours, when using PointNet vanilla as the task network.

Classification Augmentation We employed the augmentation proposed by Qi *et al.* [1]: random rotation of the point cloud along the up-axis, and jittering the position of each point by a Gaussian noise with zero mean and 0.02 standard deviation.

Reconstruction Architecture The architecture of S-NET for the reconstruction experiment follows the same basic structure as in the classification case, with 1×1 of sizes [64, 128, 128, 256, 128] and fully connected layers of sizes [256, 256, $k \times 3$]. ProgressiveNet takes the architecture of S-NET with $k = 2048$.

Sampling ratio	FPS	ProgressiveNet $(1 + \epsilon)$ matching	ProgressiveNet NN matching
1	87.3	87.3	87.3
2	85.6	85.7	85.4
4	81.2	82.4	82.3
8	68.1	76.4	78.2
16	49.4	68.3	74.4
32	29.7	50.7	61.0
64	16.3	27.8	40.0

Table 3. **Matching methods comparison.** We compare the classification accuracy of PointNet vanilla on the sampled points when using different matching methods: $(1 + \epsilon)$ and Nearest Neighbour (NN). FPS is shown for reference. NN matching is better for all sampling ratios larger than 4.

Reconstruction Optimization To train both S-NET and ProgressiveNet, we used Adam optimizer with a learning rate of 0.0005, momentum 0.9 and mini-batches of 50 shapes. The regularization weights were set to $\alpha = 0.01$, $\beta = 1$, $\gamma = 0$, $\delta = 1/64$. For ProgressiveNet, the total loss is divided by the number of sample sizes $|C_s|$. We trained the networks for 500 epochs with a Titan X Pascal GPU. S-NET training takes between 4 to 8 hours, depending on the sample size. ProgressiveNet training time is about 12 hours. No augmentation was used for reconstruction.

A.2. Matching methods

We examined two matching methods. The first one is based on nearest neighbour (NN) matching, as detailed in Section 3.2. In the second one, we adapted the implementation of the $(1 + \epsilon)$ approximation scheme for EMD [47], provided by Fan *et al.* [40]. This implementation gives a matching score for each point in one point cloud to each point in the other one. We feed the algorithm with G and P and take the points from P corresponding to the highest matching score for each point in G . In Table 3 we compare these matching methods. The accuracy is better with NN matching for large sampling ratios, and is almost the same for small ones.

NN matches every point in G with its closest point in P . This ensures that all generated points that S-NET learned are reflected in the matched set. This is not the case when using $(1 + \epsilon)$, where a generated point might be matched with a far input point, if it serves the $(1 + \epsilon)$ approximation scheme. In addition, NN matching is independent for each point, thus it guarantees to keep the order of the points produced by ProgressiveNet.

A.3. Regularization

In this section we explore the influence of the sampling regularization loss and its components. Please see Section 3.1 for the relevant background and equations.

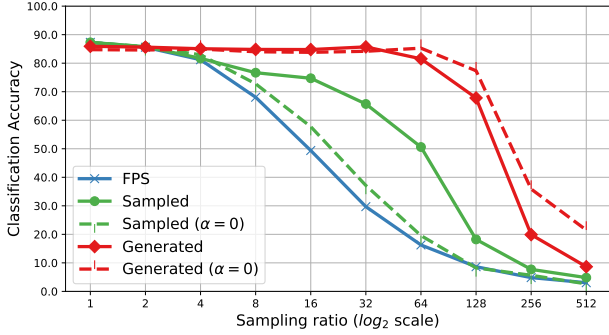


Figure 11. **Turning off the sampling regularization loss ($\alpha = 0$).** PointNet vanilla was trained on complete point clouds (1024 points) and evaluated on either S-NET’s generated or sampled points, either with or without the sampling regularization loss. Without the regularization, the accuracy using the generated points is slightly increased. However, the accuracy when using the sampled points decreases substantially.

Turning off the regularization ($\alpha = 0$) Figure 11 shows the accuracy of PointNet vanilla using either S-NET’s generated or sampled points, with and without regularization. Without the regularization, the accuracy using the generated points is slightly increased, which is expected, since the regularization sacrifices some of the classification loss to minimize the sampling loss. However, after the matching process the classification accuracy on the unregularized sampled points is much lower than when using the regularization. Without the sampling regularization, the generated points are not close to the original points and are not spread out over the whole shape. As a result, the matching process causes the sampled points to be very different from the generated points. Therefore, the sampled points do not preserve the contribution of the generated points to the task.

Turning off the linear factor for L_b component ($\delta = 0$) In Figure 12 we see the average number of unique points in the initial sampled set (after NN matching, before FPS completion) as a function of the number of generated points. For low values, the numbers are almost the same, but they diverge for higher values. Ideally, we would like to get the identity line $Y = X$, meaning that each generated point is uniquely matched with an input point without any collisions. Using a linear factor for the L_b regularization term, such that the weight is larger for large sampling sizes, improves the spread of the generated points over the shape and reduces collisions.

Turning off the L_b component ($\gamma, \delta = 0$) In Figure 13 we see the accuracy of PointNet vanilla using ProgressiveNet sampled points, either with or without the L_b regularization term. Without this term, the generated points tend to concentrate on small portions of the input point cloud, resulting in many collisions in the matching process. This

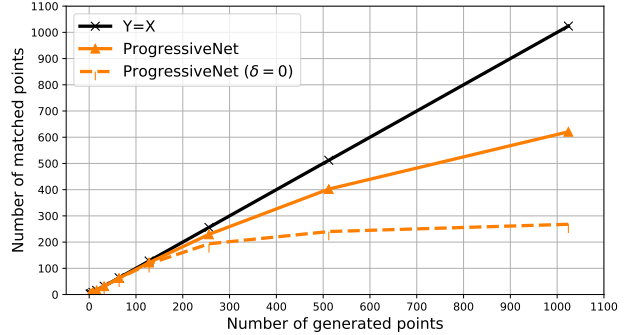


Figure 12. **Turning off the linear factor for L_b component ($\delta = 0$).** Average number of unique points in the initial sampled set (after NN matching, before FPS completion) is shown as a function of the number of generated points. Ideally, we would like to get the line $Y = X$, which means that every generated point has a unique matched input point. Setting $\delta = 0$ greatly reduce the number of unique matches.

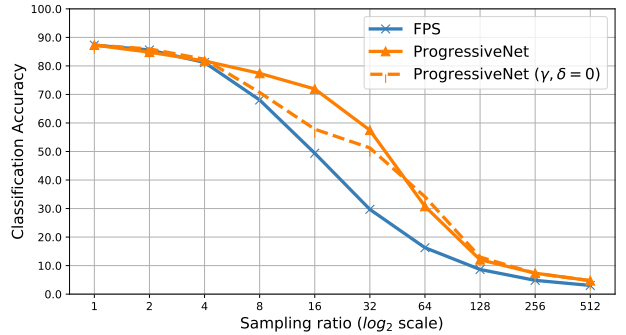


Figure 13. **Turning off the L_b component ($\gamma, \delta = 0$).** PointNet vanilla was trained on complete point clouds (1024 points) and evaluated on ProgressiveNet’s sampled points, either with or without the L_b regularization term. The results are evidently better with the L_b term.

makes the FPS completion more dominant. For low sampling ratios this is not a concern, since FPS gives decent results. For very high sampling ratios this is also not a concern, since the task loss forces the small number of generated points to spread over the input shape even without the regularization. However, for the mid-range of sampling ratios (larger than 4 and smaller than 64) the L_b regularization term is crucial. It spreads the generated points over the input point cloud, allowing for better matching and therefore better accuracy when using the sampled points.

Turning off the L_m component ($\beta = 0$) Removing the L_m term results in 3.5-fold increase in the maximal distance between the generated points and their matched points, averaged over the test set. This leads to less tight matches, resulting in up to 7% (obtained for $k = 16$) decrease in classification accuracy when using the sampled points.

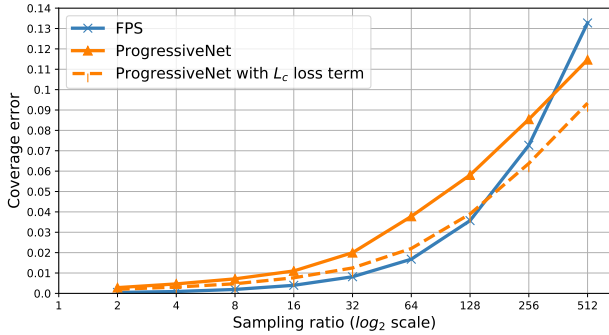


Figure 14. **Adding L_c component to the loss.** Coverage error is defined as the distance of the next FPS point. We observe that adding the L_c component closes most of the gap between FPS and ProgressiveNet.

Adding an extra regularization term FPS minimizes a geometric measure, i.e, given k points it minimizes the distance to the $k + 1$ point. We regard this measure as the coverage error. Our work found it to be an inferior proxy to minimizing the task error. Nonetheless, if in some settings the coverage error is important, we can incorporate it by adding the following loss term to Equation 5:

$$L_c(G, P) = \max_{y \in P} \min_{x \in G} \|y - x\|_2^2. \quad (9)$$

Figure 14 shows the coverage error as a function of the sampling ratio, for FPS and ProgressiveNet, as well as for ProgressiveNet when trained with the extra L_c loss term. We observe that adding the extra term removes most of the coverage error difference between FPS and ProgressiveNet.

PointNet’s accuracy when using the sampled points did not change substantially when adding the extra term, thus we conclude that it is unnecessary to minimize the coverage error in order to get a sample that is optimized for the task. However, it is not harmful to minimize this error as well, if needed.

A.4. Time and space considerations

In Section 4.1 we mentioned the influence of S-NET on inference time and memory consumption. This section elaborates on the subject. Table 4 demonstrates how we save most of the inference time by using S-NET for a small increase in memory. When fed with a complete point cloud (with 1024 points), PointNet performs 440M floating point multiplications (FLOPs). When feeding only 16 points, this number drops to 7M. S-NET that samples from 1024 to 16 points requires 35M FLOPs. Cascading them together amounts to 42M FLOPs, which is 90% reduction in inference time, compared to running PointNet on the complete points clouds. S-NET that samples 16 points requires 0.18M parameters, which is only 5% of PointNet’s memory requirement.

Network	#Parameters	FLOPs/ point cloud
PointNet-1024	3.5M	440M
PointNet-16	3.5M	7M
S-NET-16	0.18M	35M
S-NET-16 + PointNet-16	3.68M	42M

Table 4. **Time and space complexity of S-NET and PointNet.** ”M” stands for million. FLOPs stands for number of floating point multiplications. S-NET-16 stands for S-NET with output size of 16 points. S-NET-16 + PointNet-16 stands for sampling 16 points with S-NET and then classifying them with PointNet. Sampling makes for a great reduction in time complexity for a modest price in terms of space complexity.

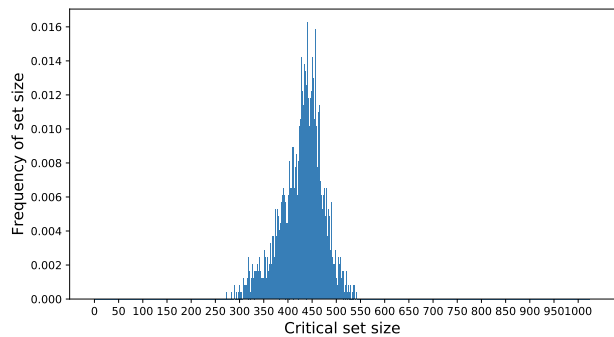


Figure 15. **Distribution of critical set sizes for PointNet vanilla over ModelNet40 test set.** The critical set size is concentrated around 429 points.

A.5. Critical set sampling

The critical set, as defined by Qi *et al.* [1], is the set of points that contributes to the max pooled feature (MPF). A proposed alternative to S-NET might be to extract the critical set and use it as the sampled point cloud. This method has three main disadvantages: first, it does not allow to control the sample size, but only sample to the size of the critical set; Second, each point cloud will be sampled to a different size, since the critical set size is not the same across the dataset, which does not allow efficient processing and storing; Third, the average critical set size of PointNet vanilla on ModelNet40 test set is 429 out of 1024 points (see Figure 15), i.e., approximately 42% of the points, which is equivalent to sampling ratio of about 2.5. On the contrary, S-NET enables control of the sample size and performs well for much larger sampling ratios. This makes the proposed alternative to S-NET not viable.

A more plausible alternative might be to sample a subset of the critical set that controls most of the features. To do so, we count the number of features that each point contributes to the max-pooled features and select the k points with most contribution. We denote this process as critical set sampling.

Sampling ratio	FPS	ProgressiveNet	Critical set sampling
1	87.3	87.3	87.3
2	85.6	85.4	87.3
4	81.2	82.3	86.3
8	68.1	78.2	76.0
16	49.4	74.4	45.6
32	29.7	61.0	22.5
64	16.3	40.0	12.5

Table 5. **Critical set sampling.** We compare the classification accuracy of PointNet vanilla for different sampling methods and sampling ratios. Critical set sampling samples the critical points that contribute to the most features to the MPF. This method is competitive for very small sampling ratios, but is not viable for larger sampling ratios.

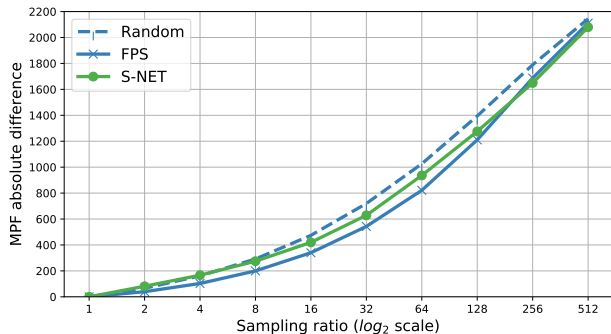


Figure 16. **Absolute difference of the max pooled feature (MPF).** Three sampling methods are compared: random, FPS and S-NET. The difference is averaged over the test set of ModelNet40. The MPF resulting from S-NET’s points differs from the original MPF like random and FPS.

In Table 5 we compare this sampling method with ProgressiveNet and FPS. We see that critical set sampling performs well for very small sampling ratios, but performs poorly for larger ratios.

Relation between our sampling and the critical set In this experiment we measured the percentage of critical set points covered by S-NET’s sampled points. We found that S-NET did not cover the critical set more than a random sample. To further investigate this issue, we computed the absolute difference between the MPF when feeding PointNet with sampled points and with the complete point cloud of 1024 points. The results are shown in Figure 16.

S-NET did not learn to sample the critical set. It also did not learn to reproduce the MPF. Instead, it learned to sample points that give better classification results, independently of the critical set. To explain this, we recall that the fully connected layers of PointNet, which process the MPF to infer the classification results, form a non-linear function. In addition, the MPF represents the shape class with redun-

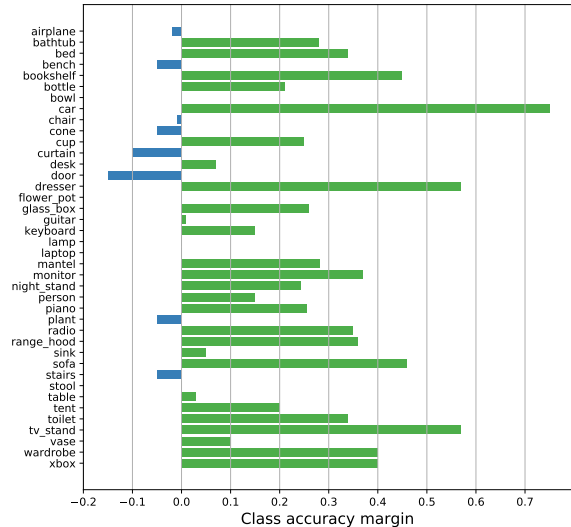


Figure 17. **Class accuracy margin.** PointNet was trained on the complete point clouds (1024 points) and evaluated on sampled point clouds of $k = 64$ points, using either FPS or S-NET sampling. S-NET achieves better results on 27 classes, with an average margin of 26%, while FPS achieves better results on 8 classes with an average margin of just 6%.

dancy (1024 floating point numbers to represent a class out of 40 classes). Thus, it is possible to find several different MPFs that results in the same classification. Therefore, there are multiple sets of points that gives similar classification results.

To further stress this point, we evaluated PointNet accuracy on the complementary set of the critical set for each shape, and the accuracy only drops by 1.5%, from 89.2% to 87.7%. We conclude that the critical set is not that critical.

A.6. Class accuracy

Up to this point we reported instance classification accuracy (also regarded as overall accuracy). Now we analyze the per class accuracy behaviour. Figure 17 shows the per-class difference in accuracy between using 64 S-NET and 64 FPS points. S-NET achieves superior results in 27 out of 40 classes (67.5%) while FPS is better in only 8 classes (20%). The results are equal for the remaining 5 classes. The average margin on the classes with superior S-NET results is 26%, compared to just 6% average margin for the classes with better FPS results. Notably, FPS achieves higher accuracy on the door class with a margin of 15% (80% vs 65%), while S-NET achieves higher accuracy on the car class with a margin of 75% (91% vs 16%).

B. Visual examples

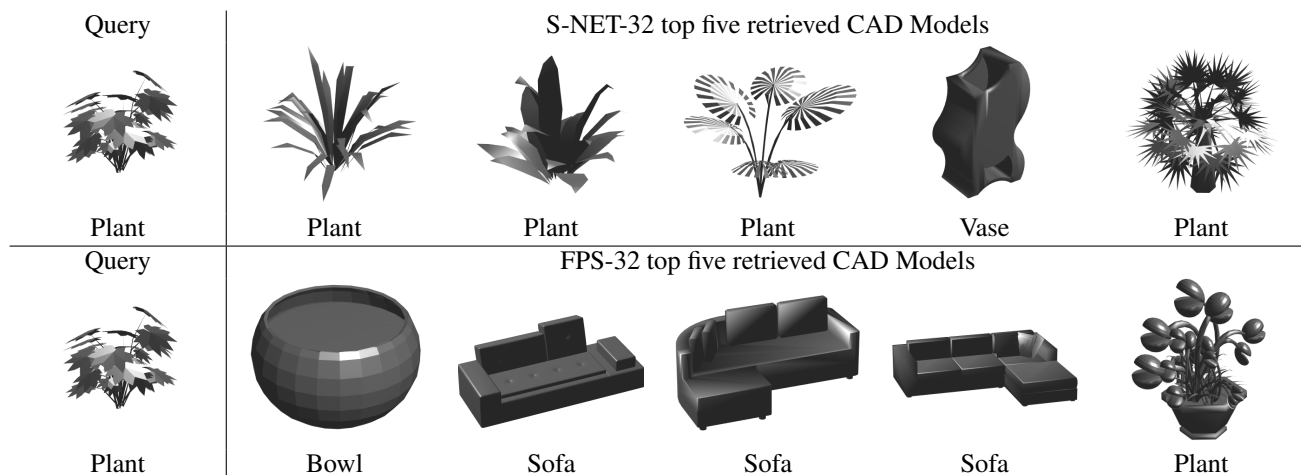


Figure 18. **Retrieval example.** We compare the top five retrievals from the test set when using 32 sampled points, either by S-NET or FPS. The sampled points were processed by PointNet (that was trained on complete point clouds of 1024 points) and its penultimate layer was used as a shape descriptor. Retrieval was done based on L_2 distance on this shape descriptor. S-NET was trained with PointNet for classification, no additional training was done for retrieval. When using S-NET, four out of five retrieved shapes are correct (plant). For FPS, only the fifth retrieved shape is correct.

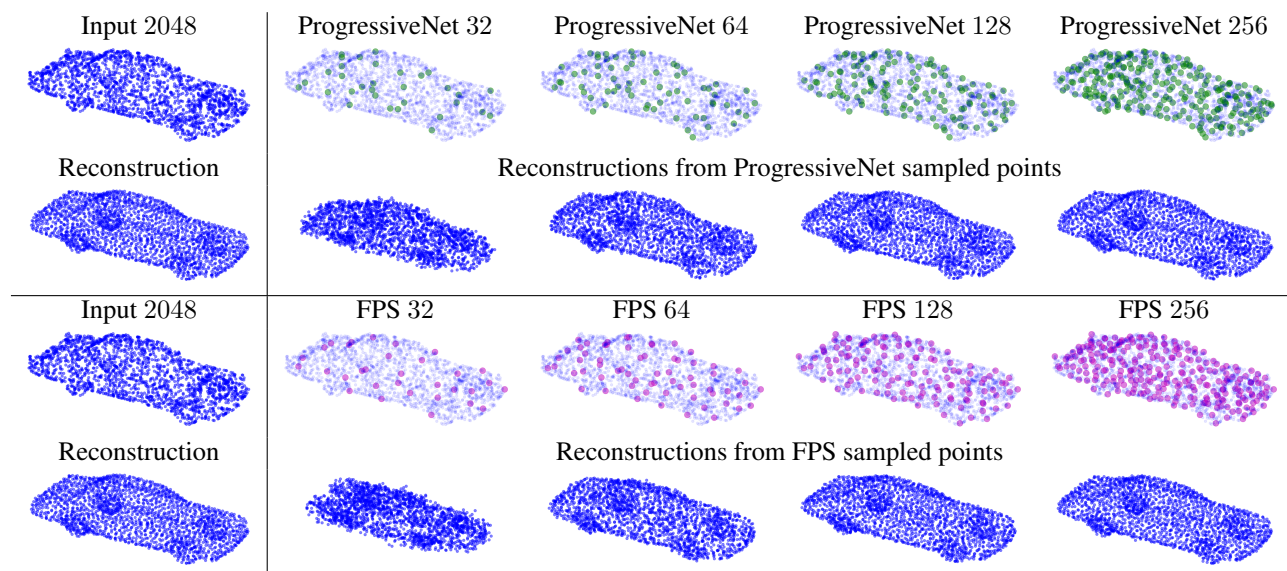


Figure 19. **Progressive sampling.** First and third rows: input point cloud and samples of different sizes by ProgressiveNet and FPS, respectively. Second and fourth rows: reconstruction from the input point cloud and from the corresponding samples. The sampled points are enlarged for visualization purpose. Even at a sample size as low as 64 points, the reconstruction from ProgressiveNet's points is visually similar to the reconstruction from the complete input point cloud. On the contrary, it takes four times more FPS points to achieve this level of similarity.

C. Extension - Progressive Autoencoder

We extended the training concept of ProgressiveNet for training a progressive autoencoder, named ProgressiveAE. That is an autoencoder whose output points are ordered according to their contribution to the reconstruction of the input point cloud.

Motivation The scheme of ProgressiveAE fits naturally to a Client-Server scenario and its benefits are three-fold: lower communication load, lower memory footprint and level-of-detail control.

Suppose that the encoder is located at the server. Instead of sending the whole point cloud (2048 3D points, which are 6144 floats) to the client, the server sends only the latent vector of 128 floats (98% reduction in communication load). The client possesses only the decoder part of ProgressiveAE, which is a progressive decoder.

According to the available memory resources and the required level-of-detail, the client may hold only the parameters corresponding to the first c output points of the decoder and reconstruct c points instead of n . For example, if the client capacity allows only reconstruction of up to 256 points, it will hold the parameters needed to calculate the first 256 output points. This amounts to more than 80% reduction in memory consumption (number of parameters) on the client side, compared to holding the complete 2048-points decoder.

Furthermore, the client can choose in real time to reconstruct an even smaller point cloud to reduce computation load. For example, calculating the first 128 output points results in almost 90% reduction in the number of floating point operations (FLOPs), compared to reconstructing the complete point cloud.

Figure 20 summarizes the time and space requirements for reconstructing point clouds of different sizes. It is much more efficient to reconstruct a point cloud to the required size than to reconstruct the complete point cloud and to then sample it to the required size.

Implementation The architecture of ProgressiveAE is the same as that of the autoencoder proposed by Achlioptas *et al.* [17], referred to as BaselineAE. The input and output of both autoencoders consist of $n = 2048$ points.

Similar to the training of ProgressiveNet, we trained ProgressiveAE with several loss terms. Each term computes the Chamfer distance between the first c point of output and the n points of the input. The overall loss is the sum of all loss terms. For this experiment we used loss terms for $C_s = \{16, 32, \dots, 2048\}$. The loss for training the BaselineAE was Chamfer distance between the input and the output point clouds of n points. All other training conditions are the same as those detailed in Section A.1.

Results To measure the reconstruction performance, we took the first c points from ProgressiveAE and computed

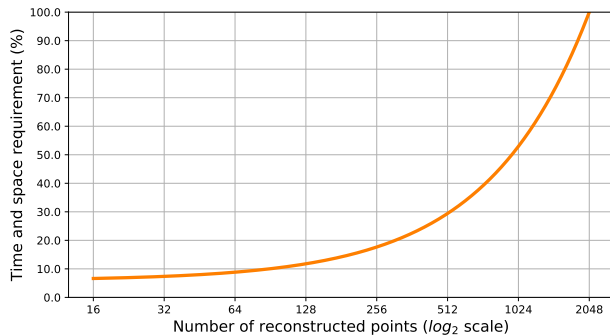


Figure 20. **Progressive decoder time and space requirements.** While a traditional autoencoder reconstructs a point cloud of a fixed size, ProgressiveAE enables real time level-of-detail management, allowing for a great reduction in inference time. In addition, the progressive decoder may hold only the parameters needed to reconstruct a point cloud smaller than the original, allowing for a reduction in memory as well.

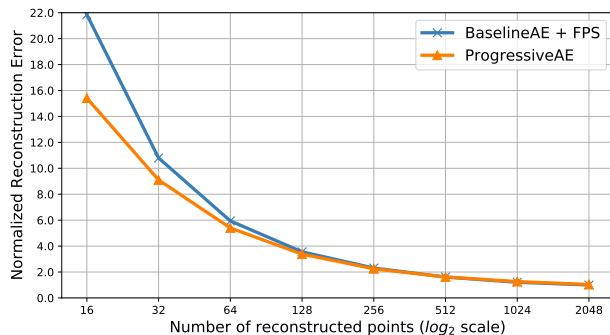


Figure 21. **Progressive autoencoder.** The Normalized reconstruction error (NRE) was computed on the test split of the data, the normalization factor is the reconstruction error when using the BaselineAE to reconstruct from complete point cloud (2048 points). ProgressiveAE has equal or better reconstruction error.

the reconstruction error as Chamfer distance from the input point cloud. As an alternative reconstruction approach with c points, we sampled c points with FPS from the output of BaselineAE. We find that ProgressiveAE has equal or better reconstruction error for for any reconstruction size (see Figure 21).

In an additional experiment, we took the encoder parameters, learned by BaselineAE, and trained only the decoder of ProgressiveAE (the variables of the layers after the max-pool operation). Interestingly, the reconstruction error for ProgressiveAE in this experiment was almost the same as in the end-to-end training of ProgressiveAE. This means that a progressive decoder can be trained to work with the encoder of an existing autoencoder, without paying any cost in terms of reconstruction quality.