

DeepPCC: Learned Lossy Point Cloud Compression

Junzhe Zhang, Gexin Liu, Junteng Zhang, Dandan Ding, and Zhan Ma

Abstract—We propose the DeepPCC, an end-to-end learning-based approach for the lossy compression of large-scale object point clouds. For both geometry and attribute components, we introduce the Multiscale Neighborhood Information Aggregation (NIA) mechanism, which applies resolution downscaling progressively (*i.e.*, dyadic downsampling of geometry and average pooling of attribute) and combines sparse convolution and local self-attention at each resolution scale for effective feature representation. Under a simple autoencoder structure, scale-wise NIA blocks are stacked as the analysis and synthesis transform in the encoder-decoder pair to best characterize spatial neighbors for accurate probability approximation of geometry occupancy and attribute intensity for point cloud compression (PCC). Experiments demonstrate that the DeepPCC remarkably outperforms state-of-the-art rules-based MPEG G-PCC and learning-based JPEG VM both quantitatively and qualitatively, evidencing that DeepPCC is a promising solution for emerging AI-based PCC. The source code will be made publicly accessible soon.

Index Terms—Point cloud compression, geometry, attribute, sparse convolution, local self-attention

I. INTRODUCTION

Large-scale point clouds have emerged as a prevalent format to realistically represent 3D objects in various applications, including AR/VR (Augmented Reality/Virtual Reality), Digital Twin, *etc.*, for an immersive experience. Typically, a large-scale point cloud sample comprises millions of unordered and sparsely distributed 3D points, with each point usually characterized by its geometry coordinate, *e.g.*, (x, y, z) in a Cartesian coordinate system, and attribute components such as RGB color, normal, and reflectance. For a color point cloud frame that has 1 million points and uses 8-bit precision to represent both geometry and RGB attributes of each point, it requires about 6 MB data size per frame, which apparently cannot be sustained reliably for a variety of access networks used in practices [2]. This thus urges the development of high-efficiency Point Cloud Compression (PCC) solution for uncompromised service provisioning [3].

A. Background and Motivation

The lossy compression of point clouds is a classical rate-distortion (R-D) optimization problem in information theory [4], where the goal is to obtain a compact representation of the input sample with the least bit rate consumption under a given distortion metric, *e.g.*, MSE (Mean Square Error) or MS-SSIM (Multiscale Structural Similarity [5]). In principle, the derivation of this compact representation of an input point cloud is mainly determined by how to properly organize its

J. Zhang, G. Liu, J. Zhang, and D. Ding are with the School of Information Science and Technology, Hangzhou Normal University, Hangzhou, Zhejiang 311121, China. Email: DandanDing@hznu.edu.cn

Z. Ma is with the School of Electronic Science and Engineering, Nanjing University, Nanjing, Jiangsu, 210093 China. Email: mazhan@nju.edu.cn

data to sufficiently exploit spatial (and/or temporal) correlations for redundancy reduction.

Given the unordered and sparse nature of point clouds, an enormous amount of effort has been paid over the past two decades to effectively organize spatial neighborhoods for correlation exploration. Such neighborhoods could be formed through the use of 1D transversal, 2D projection, or 3D space indexing. A number of pioneering explorations have been developed, as reviewed in [2], leading to the establishment of two international standards for PCC under the MPEG committee in 2017: Video-based PCC (V-PCC) and Geometry-based PCC (G-PCC) [3], [6], [7].

V-PCC first projects a 3D point cloud frame onto a set of perpendicular 2D planes [8], and then employs conventional video coding standards [9], [10] to process the resulting sequence of projected 2D images. In contrast, G-PCC separates the compression of geometry and attribute information, using a 3D octree model or a trisoup model to encode the geometry information and employing Region-Adaptive Hierarchical Transform (RAHT) [11] or hierarchical nearest neighbor prediction-based Lifting Transform (Predlift) to represent the attributes (such as RGB colors) on top of the restored geometry.

Owing to the advances of deep learning technologies, recent years have witnessed the exponential growth of learning-based PCC solutions [12]. More importantly, both MPEG and JPEG committees have launched explorations of AI (Artificial Intelligence) based PCC, with the aim of establishing it as the next-generation standard [13], [14]. It is worth pointing out that the JPEG committee has released the very first learning-based Verification Model (VM), leading the standardization development of AI-based PCC for the future¹.

This work proposes DeepPCC, a new learning-based lossy PCC approach for the compression of large-scale object point clouds.

B. Our Approach and Contribution

Unlike existing works that focus on either geometry or attribute information compression [12], the proposed DeepPCC can process both geometry and attribute components of an input point cloud. Following the conventions used in MPEG G-PCC standard [15], we first compress the geometry, and then the attribute components such as colors, assuming full knowledge of the corresponding geometry. Voxelized point clouds are utilized in DeepPCC, which is consistent with the processing of MPEG G-PCC and JPEG VM. However, the point-based compression method without voxelizing raw points is beyond the scope of this work; interested readers can refer to [16]–[21] for more details.

¹<https://jpeg.org/jpegplen/poincloud.html>



Fig. 1: **Visual Reconstructions** of point clouds compressed using DeepPCC, MPEG G-PCC, and JPEG VM. Both geometry and color attributes are lossy compressed. Two G-PCC models, e.g., TMC13v6 and TMC13v19, and the latest JPEG VM are used. Bitrates are made as close as possible to have a fair comparison. However, as these methods do not have rate control yet, bitrate mismatch happens inevitably. Yet, the proposed DeepPCC shows clear gains for both quantitative and qualitative measures. bpp - bits per point; PSNR: Peak-to-signal-noise ratio; D1: point-to-point PSNR; Y: Luma PSNR; PCQM: a full-reference quality metric for colored point cloud [1]. Please zoom in for more details.

In theory, the efficiency of PCC highly depends on the effectiveness of exploiting inter-point or inter-voxel correlations², which is challenging due to the unordered and sparse nature of 3D points. To address this challenge, we first assume the use of a uniform voxel representation model for 3D position indexing [22]. For a given point, we can easily access its neighbors in a $n \times n \times n$ cube centered at it. Straightforwardly, dense 3D convolutions can be applied on each point to aggregate and embed spatial neighbors for compact representation, as studied in [23], [24], which however introduces unbearable space and time complexity due to the uniform computations on every voxel regardless of its occupancy state. Although block partition or patch slicing methods can be utilized to reduce the complexity to some extent, their performance is often constrained to local optima [25].

Instead, we propose to apply sparse convolution [26] and local self-attention [27] to strictly limit the computations only to positively occupied voxels (POV) in close proximity, which reduces complexity significantly and enables us to process the entire large-scale point cloud frame directly. The proposed local self-attention performs computations only among K nearest POV neighbors (KNN) of a specific POV, which thus is referred to as the KNN self-attention for convenience.

In DeepPCC, we take advantage of sparse convolution

and KNN self-attention by either directly stacking them or allocating them in parallel branches to achieve efficient Neighborhood Information Aggregation (NIA). Since the POVs in a point cloud frame are sparsely distributed, it is highly likely that only a limited number of POVs will fall within the fixed-size receptive field of sparse convolution, which may not provide sufficient information to accurately characterize neighborhood correlations for compression. In such cases, KNN self-attention can gather K neighboring POVs without being limited by the fixed-size receptive field, providing the necessary information for NIA. Yet, using a smaller K may not capture sufficient POVs nearby, while using a larger K could lead to a sharp increase in computational complexity. Therefore, to balance the complexity and efficiency, we often set a reasonable and fixed value of K . To further enhance information embedding capacity, we propose to progressively downscale the point cloud (e.g., 3 times in this work), and at each resolution scale, the NIA block is incorporated.

As illustrated in Fig. 2, the DeepPCC adopts a simple autoencoder structure, in which it stacks scale-wise NIA blocks to form respective analysis and synthesis transform in paired encoder and decoder. To help the audience understand the DeepPCC better, we draw two parallel branches: one branch for the progressive downscaling³ to generate multi-

²After the voxelization, raw points are mapped to positively-occupied voxels (POV). For simplicity, we may interchangeably use the term “point” and “POV”.

³As will be detailed in subsequent sections, dyadic downsampling at each axis in Cartesian coordinate system is applied for the geometry component, while average pooling is used for corresponding attributes.

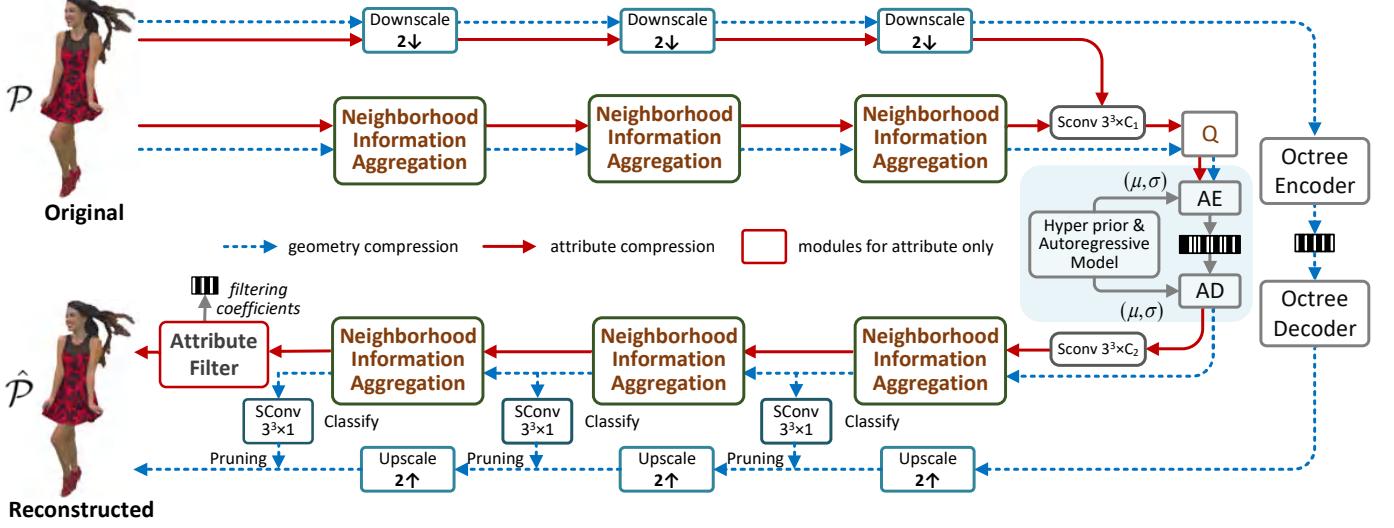


Fig. 2: **DeepPCC** used for lossy point cloud compression. For both geometry compression (indicated by the blue dotted lines) and attribute compression (indicated by the red solid lines), multiscale NIA blocks are devised to best characterize neighborhood correlation for geometry- and attribute-related latent features in encoding. In decoding, geometry occupancy is scale-wisely decoded and refined using lower-scale information, while the attribute features are decoded to the final resolution in one shot. An Attribute Filter is additionally introduced to enhance the quality of attribute reconstruction. SConv $k^3 \times c$: Sparse convolution with $k \times k \times k$ kernel and c channels output; $s \downarrow/s \uparrow$: downscaling/upscaling with a ratio of s at $x-$, $y-$, and z -axis; (μ, σ) : learned entropy parameters.

resolution point clouds; and the other branch for applying the corresponding NIA in feature domain at each resolution scale.

In lossy geometry encoding, the geometry coordinates of POVs at the lowest scale (after downsampling) are losslessly coded using the G-PCC coder for simplicity, while the corresponding latent features at the bottleneck are quantized and encoded using adaptive contexts. In the lossy attribute encoding, the latent features are also encoded using adaptive contexts, given losslessly or lossy compressed geometry priors. During the decoding process, the geometry occupancy is gradually refined from the lower scale to the higher scale. Such a scale-wise geometry reconstruction and refinement not only provides better quality but also significantly reduces the space and time complexity [25]. On the other hand, the attribute features are decoded to the final resolution in one shot. To alleviate the compression noise, an Attribute Filter is applied at the end to enhance the reconstruction quality.

Currently, the proposed DeepPCC supports three modes: DeepPCGC for lossy geometry compression only, DeepPCAC for lossy attribute compression with losslessly compressed geometry, and DeepPCC for lossy compression of both geometry and attribute under the same framework shown in Fig. 2.

The main contributions of this paper are summarized below:

- The proposed DeepPCC is a simple yet efficient learning-based lossy PCC approach that can process both geometry and attribute components under a unified framework, offering a competitive alternative solution to existing rules-based MPEG G-PCC and learning-based JPEG VM solutions;
- DeepPCC demonstrates significant performance gains over both MPEG G-PCC and JPEG VM across a va-

TABLE I: Notations

Abbreviation	Description
PCC	Point Cloud Compression
PCG	Point Cloud Geometry
PCA	Point Cloud Attribute
POV	Positively-Occupied Voxel
NIA	Neighborhood Information Aggregation
KNN	K Nearest Neighbors
RAHT	Region-Adaptive Hierarchical Transform
PSNR	Peak-to-signal-noise ratio
BD-BR	Bjøntegaard Delta Rate
BD-PSNR	Bjøntegaard Delta PSNR

riety of testing conditions. For the lossy compression of both geometry and attribute information, DeepPCC offers more than 60% BD-BR gains in comparison to the G-PCC (TMC13v19) when measuring the overall distortion using PCQM (see Table VI); Similar gains are reported when using the JPEG VM as the anchor. Qualitatively, DeepPCC produces more visually-pleasing reconstructions (see Fig. 1).

- The superior performance of DeepPCC comes from the utilization of multiscale neighborhood information aggregation, which enables the efficient characterization of spatial correlations for compact representation. Moreover, using sparse convolution and KNN self-attention limits the computation of DeepPCC to positively-occupied voxels only, making it attractive to practitioners due to its affordable complexity.

Table I gives the notations and abbreviations frequently used in this work.

II. RELATED WORK

This section briefly reviews related publications that focus on the use of learning tools to improve the efficiency of PCC. In the meantime, a recent survey paper [12] on learning-based algorithms proposed in the past years serves as a good reference. Additionally, an article by Cao et al. in [2] can help audiences understand classical rules-based algorithms in PCC.

A. Point Cloud Geometry Compression (PCGC)

Many works were devoted to studying point cloud geometry compression (PCGC). To exploit the correlation across points, an appropriate presentation model is first required to construct the local neighborhood for information characterization and embedding. There are three popular representation models: point-based model [16], [18]–[20], [28]–[30], uniform voxel representation [23], [24], [31], octree scheme [32]–[36], and multiscale sparse representation [22], [25].

These approaches utilized the deep neural networks (DNNs) involving various computational layers, such as 3D dense convolution, sparse convolution, self-attention, or MLP (MultiLayer Perceptron), to predict the geometry occupancy probability of each voxel or point through the aggregation of neighborhood priors. As seen, their effective spatial neighborhood is defined by the representation models aforementioned. Among them, the work proposed by Guarda *et al.* [24] was now part of the JPEG Pleno PCC software model (*i.e.*, JPEG VM) under development; the SparsePCGC developed by Wang *et al.* [22] is under investigation in MPEG 3D Graphics group as a potential solution for next-generation AI PCC. Both approaches have shown significant performance improvement over the G-PCC in terms of point cloud geometry compression.

Instead of developing end-to-end solutions, another avenue is to restore the reconstruction quality by compensating for G-PCC-induced compression impairments. For example, GRASP-Net [37], DGPP [38], and GRNet [39] were proposed to remove artifacts induced by the underlying G-PCC coder.

B. Point Cloud Attribute Compression (PCAC)

Learning-based point cloud attribute compression (PCAC) just attracted attention recently. Compared with the geometry occupancy that presents binary status (occupied or not), attribute components (*e.g.*, RGB colors) exhibit larger variations across neighboring points, which, to some extent, makes them more difficult to characterize and embed in the raw data format. To this end, existing works often applied transforms (*e.g.*, Graph Fourier Transform [40], [41] and RAHT [11]) upon attributes to better represent signals in the transform domain. Note that the RAHT which is basically an adaptive Haar wavelet transform is adopted in G-PCC.

Some methods directly use stacked DNNs as the transform in an end-to-end framework. Two typical examples are the Deep-PCAC [42] and SparsePCAC [43]. In Deep-PCAC, MLPs were used to directly process unorganized, unordered points, with each input point cloud divided into 2048-point patches for processing. As seen, Deep-PCAC is constrained to exploiting point correlations within local patches and is prone

to falling into local optima. Consequently, their results suffered a significant loss in BD-PSNR (up to 5.67 dB) compared to an earlier G-PCC version (TMC13v10). SparsePCAC used a simple autoencoder structure with sparse convolutions, offering better performance and requiring much less complexity than Deep-PCAC. Although SparsePCAC still presented inferior performance to the most recent G-PCC version (TMC13v19), it shows a promising direction for attribute compression using learning-based methods.

Instead of solely relying on the DNNs to characterize neighborhood correlations, both Nguyen *et al.* [35] and Wang *et al.* [44] proposed to use local neighbors explicitly structured by underlying representation models, such as octree or multiscale sparse tensor, to predict attribute probability. Their methods demonstrated superior performance compared to G-PCC in lossless mode. However, the decoding runtime of [35] is unbearable due to the use of autoregressive contexts in the entropy coder, while the decoding time of [44] is comparable to that of G-PCC because of its multistage parallelism strategy.

Similar to explorations in PCGC, some works employed DNNs to enhance the G-PCC attribute compression. For example, 3DAC [45] used an MLP-based entropy model to better encode the RAHT coefficients of G-PCC, which resulted in some gains on samples from ScanNet and Semantic KITTI datasets. Another example is CARNet [46] that developed a neural filter on top of the G-PCC for artifact removal and reported more than 20% compression gains.

Remark. As seen, most works focused on compressing either geometry or attribute information individually. Although Nguyen *et al.* [35] provided a solution to encode geometry and attribute under the same framework, it was limited to the lossless mode and therefore unsuitable for applications that desire lossy compression. Moreover, its extremely-long decoding complexity makes it impractical for real use.

III. DEEPPCC

The journey of DeepPCC started from PCGCv2 [25], PCGFormer [27], and recent SparsePCAC [43]. The former two only discussed the lossy geometry compression while the latter focused on the attribute coding. Although PCGCv2/PCGFormer showed impressive performance improvements to the octree-based geometry compression in G-PCC, SparsePCAC was still largely inferior to the RAHT-based attribute coding of G-PCC. This work proposes a unified framework for lossy geometry and attribute compression, providing a comprehensive and promising solution for next-generation AI-based PCC. Fig 2 sketches the general structure of DeepPCC. As DeepPCC separates the processing of geometry and attribute coding, we next present more details on each.

A. Lossy DeepPCGC

The lossy geometry compression of DeepPCC, *a.k.a.*, DeepPCGC, is illustrated using the blue dotted lines in Fig. 2. Network implementation with detailed layers of DeepPCGC is given in the supplementary material.

Given an input point cloud frame \mathcal{P} , we use $\{C_X, F_X\}$ to represent its geometry part PCG, where C_X collects the

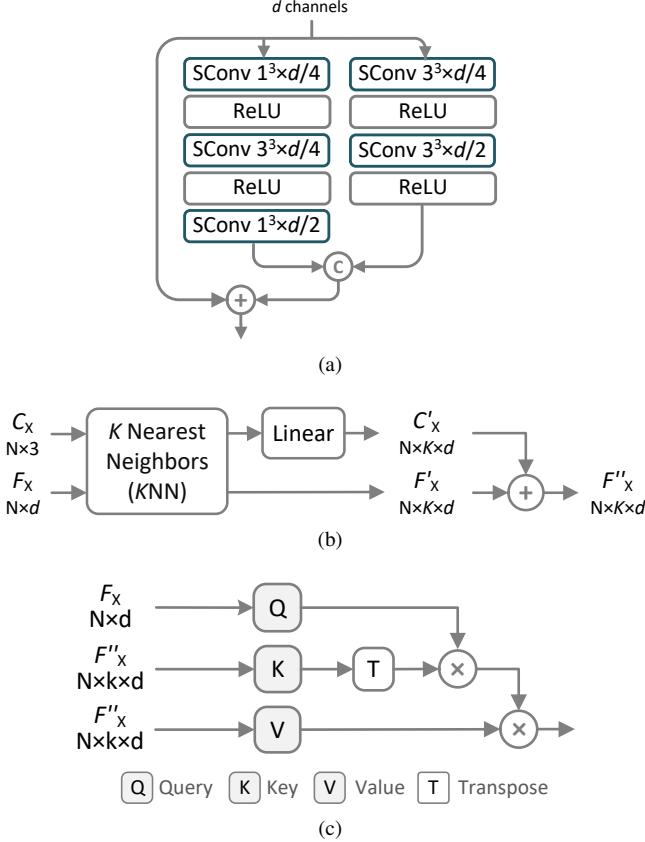


Fig. 3: NIA Modules. (a) **IRN** stacks sparse convolution (SConv) layers for information embedding. (b) **KNN** constructs adaptive local neighborhood by searching K nearest neighbors of each POV. (c) **KNN Self-Attention** used for local information aggregation where an attention map is generated for every group of K POVs in close proximity. Attentive features are then produced by multiplying the input value with the corresponding attention map.

coordinates of POVs and F_X is the feature vector associated with POVs. Specifically, for the input PCG, all elements in F_X are set to “1” (or TRUE) for the indication of positive occupation. F_X will be adapted by adjusting the channels from one scale to another to embed information aggregated from NIA blocks.

Encoder Transform. For the encoding, input geometry coordinates are dyadically downsampled to generate multi-resolution PCGs. At each scale, the NIA block $G_{\text{NIA}}(\cdot)$ is comprised by stacking an Inception-Residual Network (IRN) layer (see Fig. 3a), a KNN Self-Attention layer (see Fig. 3b and 3c), multiple sparse convolutional (SConv) layers, and ReLU layers for neighborhood correlation characterization and feature embedding. In this work, spatial resolution is downsampled three times to produce the thumbnail PCG C_Y , which is implemented using Sconv with a stride of 2. Correspondingly, three NIA blocks are devised to finally produce the latent feature F_Y for entropy coding, *i.e.*,

$$\{C_Y, F_Y\} = G_{\text{NIA}}(G_{\text{NIA}}(G_{\text{NIA}}(\{C_X, F_X\}))). \quad (1)$$

Note that when processing features in NIA blocks, correspond-

ing geometry coordinates are implicitly inferred.

Entropy Coding. At the bottleneck, the thumbnail PCG C_Y is losslessly encoded using the G-PCC (octree) codec; while the corresponding F_Y are quantized and entropy-coded using adaptive contexts learned from the data. Currently, we use joint hyperpriors and autoregressive neighbors in context modeling for optimal performance at the cost of the increased complexity [47]. We can simplify the context model by just using the hyperpriors or even fully factorized density model [48] for much faster processing speed. Following the Laplacian distribution used in [47], [49], entropy probability approximation can be translated as the derivation of entropy parameters (μ, σ) .

Decoder Transform. For the decoding, it basically mirrors the encoder processing to progressively upscale, classify, and reconstruct PCGs. Specifically, at each scale, the NIA is reused to generate the occupancy probability for each voxel at the current scale. A classifying and pruning operation follows to sort the top M points as restored coordinates. Note that M is a preset point number counted for each scale in the encoder as ground truth and encapsulated in the bitstream for decoding. As seen, lower-scale geometry reconstruction serves as the conditional prior for such scale-wise refinement, which not only improves the reconstruction quality with less distortion but also significantly reduces the space and time complexity [25].

Loss Function. Recalling that geometry occupancy is a binary state, a hierarchical binary cross-entropy (BCE) loss is then applied to train the DeepPCGC model, *i.e.*,

$$\mathcal{L} = R + \lambda_g \cdot D = R_{\hat{y}} + R_{\hat{z}} + \lambda_g \cdot \mathcal{L}_{\text{BCE}}, \quad (2)$$

where λ_g controls the R-D trade-off. $R_{\hat{y}}$ and $R_{\hat{z}}$ are the bitrate of latent features \hat{y} and hyperprior encoder \hat{z} in DeepPCGC, respectively. The loss \mathcal{L}_{BCE} is defined as

$$\mathcal{L}_{\text{BCE}} = \sum_i -(x_i \log_2 p(x_i) + (1 - x_i) \log_2(1 - p(x_i))), \quad (3)$$

where x_i denotes the true occupancy state of a voxel and $p(x_i)$ is the estimated occupancy probability of this voxel. On the other hand, $1 - x_i$ represents the non-occupied voxel with the corresponding probability $1 - p(x_i)$. Such a BCE loss is also widely used in other learning-based point cloud geometry compression works [12].

B. Lossy DeepPCAC

The lossy attribute compression of DeepPCC, *a.k.a.*, DeepPCAC, is illustrated using red lines in Fig. 2, which shares a similar encoder-decoder structure (*e.g.*, transform and entropy coder) with DeepPCGC. Its Detailed network structure is also provided in the supplementary material. As aforementioned, geometry is losslessly compressed in this mode.

Encoder Transform. At the encoder, **NIA** blocks are also applied scale-wisely to effectively characterize and embed attribute information (*e.g.*, RGB colors in this study) of neighboring spatial regions for the derivation of attribute-related latent features. The same equation in Eq. (1) can be used to formulate the attribute information aggregation, where F_X

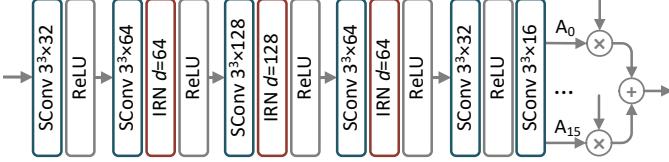


Fig. 4: **Attribute Filter** that uses filtering coefficients to guide the neural network for adaptive filtering. As such, the appearance of the reconstructed point cloud will be adjusted according to the input, uncompressed point cloud. The filtering coefficients will be coded and transmitted in the bitstream. 16 coefficients are employed in this work.

and F_Y represent the original attribute tensor and its relevant latent features at the bottleneck, respectively. The geometry coordinates are implicitly inferred during the processing steps.

Different from DeepPCGC which sequentially stacks of IRN and KNN Self-Attention layers for geometry compression, DeepPCAC allocates them in parallel branches to separately characterize attribute correlations in the local spatial neighborhood. The latent features obtained from each processing branch are then concatenated with the thumbnail point cloud downsampled three times at the bottleneck for entropy coding. Note that the attribute downscaling is achieved through average pooling, *i.e.*, eight voxels are squeezed to one voxel, with its intensity assigned using the average attribute value from the corresponding eight voxels.

Entropy Coding The same entropy engine is shared with DeepPCGC and DeepPCAC models. However, their context probabilities generated are different and substantially rely on their respective data.

Decoder Transform. At the decoder, symmetric NIA blocks are applied to process attribute-related features scale-wisely for the restoration of attribute intensities. Similarly, two parallel branches are devised to respectively apply the IRN and KNN Self-Attention layers. These two branches remain separate until they reach the last scale, where they are merged to produce the attribute reconstruction. Compression artifacts are inevitably introduced, and thus an additional Attribute Filtering is augmented next to further improve the reconstruction quality.

Attribute Filtering. We devise a network using stacked sparse convolutions (Fig. 4) to generate 16 reconstruction candidates $\mathbf{R} = \{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{15}\}$, which are then linearly combined by a set of filtering coefficients $A = \{a_0, a_1, \dots, a_{15}\}$ for the final reconstruction $\hat{P} = \mathbf{R}A^T = a_0\mathbf{r}_0 + a_1\mathbf{r}_1 + \dots + a_{15}\mathbf{r}_{15}$. These filtering coefficients are scaled to integers, compressed using the fixed-length codes, and explicitly signaled in the bitstream. The derivation of A can refer to our previous work [50].

Loss Functions. DeepPCAC is trained using a two-stage training approach. In the first stage, we use the Lagrangian loss for R-D optimization to train the multiscale network (without Attribute Filter) only:

$$\mathcal{L}_{s_1} = R + \lambda_a \cdot D = R_{\hat{y}} + R_{\hat{z}} + \lambda_a \cdot \left\| P - \tilde{P} \right\|_2^2, \quad (4)$$

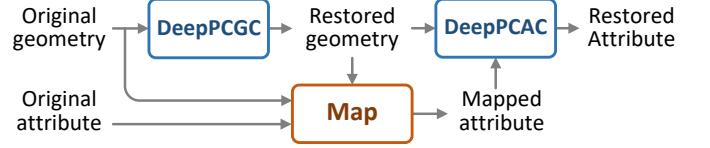


Fig. 5: **DeepPCC for Unified Lossy Coding** where DeepPCGC and DeepPCAC are connected consecutively. we first use DeepPCGC to encode and restore geometry coordinates and then map the original attributes to these restored coordinates for attribute coding using DeepPCAC.

where λ_a is the parameter that controls the R-D trade-off. $R_{\hat{y}}$ and $R_{\hat{z}}$ are the bitrate of latent features \hat{y} and hyperprior encoder \hat{z} in DeepPCAC, respectively. The MSE distortion is derived between ground truth $P \in \mathcal{P}$ and reconstructed output $\hat{P} \in \tilde{\mathcal{P}}$ in the YUV color space (without filtering).

Secondly, we fix the weighting parameters of the multiscale model and train the Attribute Filter only:

$$\mathcal{L}_{s_2} = \left\| P - \hat{P} \right\|_2^2 = \left\| P - (a_0\mathbf{r}_0 + \dots + a_{15}\mathbf{r}_{15}) \right\|_2^2, \quad (5)$$

with $A = \{a_0, \dots, a_{15}\}$ obtained by the least square optimization and $\{\mathbf{r}_0, \dots, \mathbf{r}_{15}\}$ generated by the network of Attribute Filter (see Fig. 4).

C. Unified Lossy DeepPCC

As for the lossy compression of both geometry and attribute components of a given point cloud, DeepPCC connects pre-trained DeepPCGC and DeepPCAC models, as illustrated in Fig. 5. More specifically, restored geometry coordinates from DeepPCGC are passed to DeepPCAC as implicit geometry prior, with which original attributes are mapped onto these lossy-compressed coordinates for subsequent processing.

To perform the mapping, DeepPCC utilizes a KNN-based operation. For each reconstructed coordinate, DeepPCC identifies its K_M nearest points in the original, uncompressed point cloud and computes the average attribute value of these K_M points as its attribute. In our implementation, we set $K_M = 3$. Using a larger K_M may lead to oversmoothed attribute intensity, which is not recommended.

IV. EXPERIMENTS

A. Testing and Training Conditions

Datasets. We randomly sampled 10,000 models from ShapeNet dataset [51] to first build our training dataset for point cloud geometry compression. This dataset is referred to as the ShapeNet PCG training dataset, in which we uniformly set 8-bit geometry precision. Such a collection of ShapeNet PCG samples is used to train the geometry compression model of DeepPCC, *i.e.*, DeepPCGC.

To establish the training dataset for the attribute compression of DeepPCC (*i.e.*, DeepPCAC), we followed the methodology discussed in SparsePCAC [43], in which color images randomly selected from COCO dataset [52] were projected onto the ShapeNet PCG samples to synthesize the color attributes.

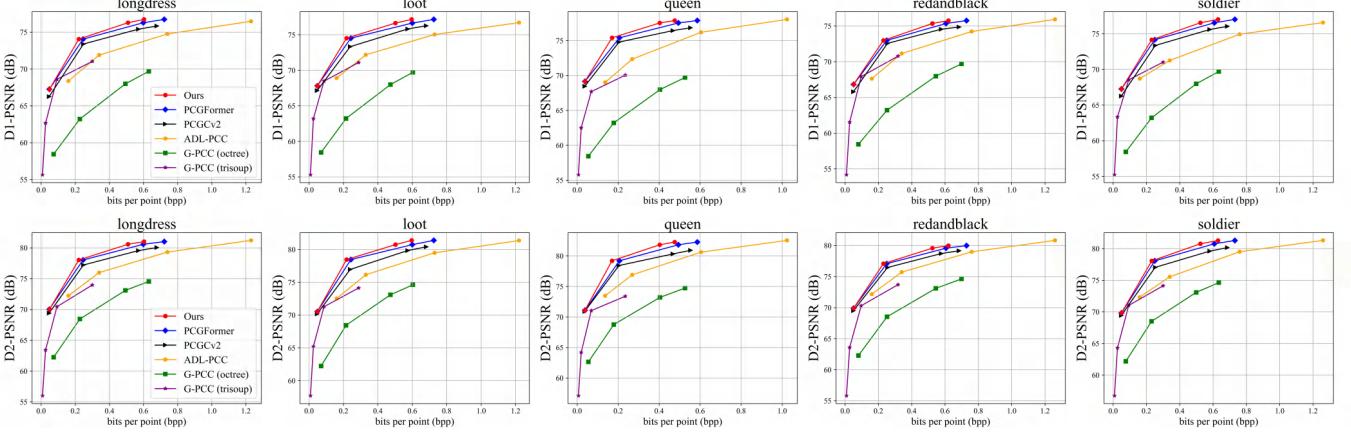


Fig. 7: **R-D Comparison across Different Geometry Compression Methods.** All neural models are trained under the same conditions. R-D curves measured in D1 and D2 PSNR are both provided.

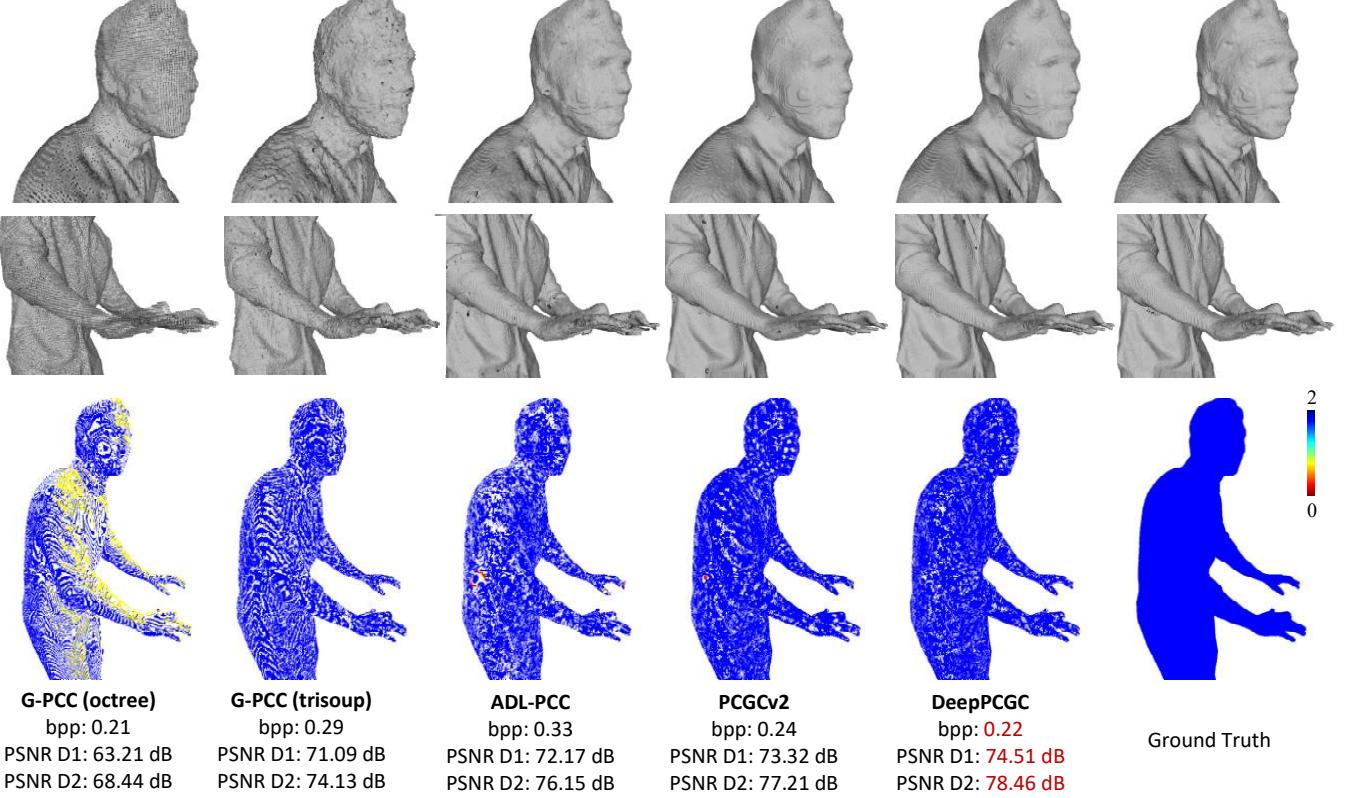


Fig. 8: **Visualized Reconstructions** of samples compressed using DeepPCGC, G-PCC (octree), G-PCC (trisoup), ADL-PCC [24], and PCGCv2 [25]. Please zoom in for details.

- **DeepPCAC** for lossy attribute coding but assuming the lossless compression of the corresponding geometry part;
- **DeepPCC** for lossy compression of both geometry and attribute components.

B. DeepPCGC Performance

BD-BR Comparison. Table II presents the BD-BR gains of DeepPCGC against state-of-the-art methods, including G-PCC (octree), G-PCC (trisoup), ADL-PCC [24], PCGCv2 [25], and PCGFormer [27]. The latest G-PCC, *i.e.*, TMC13v19, is

employed. As seen, compared with the anchor G-PCC (octree), DeepPCC achieves 91.31% (87.74%) BD-BR reduction in D1 (D2) measurement on average. Moreover, DeepPCC significantly outperforms other learning-based methods, *e.g.*, it surpasses the ADL-PCC by 68.16% (63.34%) BD-BR gains. Note that ADL-PCC was adopted in JPEG VM for geometry coding. R-D curves plotted in Fig. 7 consistently show evidence of the leading performance of DeepPCGC across a variety of bitrates.

The DeepPCGC achieves a BD-BR reduction of 26.91%

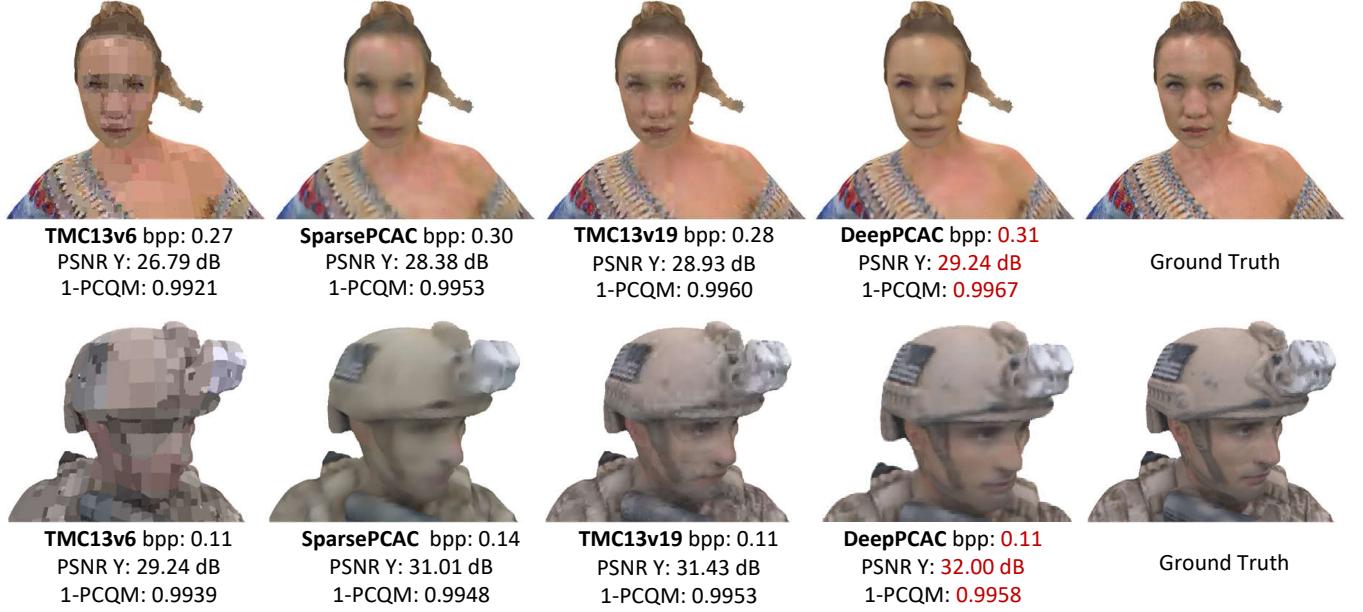


Fig. 10: **Visualized Reconstructions** of samples compressed using DeepPCAC, TMC13v6, SparsePCAC, and the latest TMC13v19. Please zoom in for details.

TABLE IV: PCQM Performance of DeepPCAC against State-of-the-art Methods

Class	Point Cloud	TMC13v6	SparsePCAC	TMC13v19
		BD-BR (%)	BD-BR (%)	BD-BR (%)
MPEG 8IVFB	longdress	-49.68	-27.27	-19.20
	loot	-10.95	-23.30	34.15
	queen	-51.03	-55.29	-34.22
	redandblack	-54.56	-38.01	-22.38
MPEG Owlii	soldier	-47.86	-35.69	-15.12
	basketball	-48.33	-36.41	-9.87
	dancer	-56.26	-45.50	-24.55
	exercise	-33.23	-27.51	-24.53
	model	-50.37	-32.48	-14.42
	Average	-44.70	-35.72	-9.01

D curves are exemplified in Fig. 9 to supplement the average BD-BR/BD-PSNR measures.

For the YUV space, 0.72% BD-BR loss is observed due to the substantial loss in *loot* sequence. This is potentially caused by the irregular texture of the shirt, which can not be well characterized by the current DeepPCAC model trained from a vastly different synthesized ShapeNet dataset. An interesting problem is to generate training samples that are more realistic and closer to the point clouds used in practice. To validate our hypothesis, Table III also reports the performance using the synthesized ShapeNet test dataset. As seen, the DeepPCC uniformly outperforms state-of-the-art methods by a large margin, *e.g.*, it gains as large as 32.14% and 34.99% BD-BR over the TMC13v19 and 19.05% and 21.52% over the SparsePCAC in Y and YUV spaces, respectively.

Furthermore, DeepPCAC also exhibits superiority when evaluated by the PCQM metric [1]. The PCQM is a linear combination of geometry and color features, which are computed between a distorted point cloud and its corresponding reference. As presented in Fig. 9 and Table IV, on average, DeepPCAC provides significant performance gains.

TABLE V: Experimental Settings in G-PCC and DeepPCC

GPCC	Octree (PQS*)	0.25	0.5	0.75	0.875
	RAHT (QP**) 35	35	35	36	36
DeepPCC	DeepPCGC (λ_g)	0.1	1	5	10

* PQS-Position Quantization Scale; **QP-Quantization Parameter.

In both Table III and Table IV, almost 40% BD-BR gains are reported when comparing DeepPCAC with SparsePCAC, which is mainly due to the use of proposed NIA module and Attribute Filter.

Qualitative Comparison. Furthermore, we illustrate the qualitative results of different methods in Fig. 10. Clearly, the visualization of SparsePCAC looks very blurry and exhibits some color casting (*e.g.*, on the lady's face and feet). The reconstruction of G-PCC TMC13v19 contains severe blocking artifacts, whereas DeepPCAC provides a very smooth and clear reconstruction, which is visually more appealing than the other methods.

D. DeepPCC Performance

Overall Performance. We form a unified DeepPCC solution for the lossy compression of both geometry and attribute by connecting the DeepPCAC and DeepPCGC models. We further conduct a comprehensive study on the compression performance of DeepPCC, compared with the latest G-PCC (TMC13v19) which is also comprised of two modules, namely G-PCC (octree) and G-PCC (RAHT). To ensure a fair comparison, we try our best to set comparable geometry bit rates in DeepPCGC and G-PCC (octree), and also similar attribute bit rates in DeepPCAC and G-PCC (RAHT), which eventually leads to comparable total bit rates. Table V details the encoder settings in DeepPCC and G-PCC. However, it is important to note that such settings may not be optimal for these two

TABLE VI: Compression Gains of DeepPCC to G-PCC (TMC13v19)

Point Cloud	DeepPCGC vs. TMC13v19				DeepPCAC vs. TMC13v19				DeepPCC
	D1		D2		Y		YUV		PCQM BD-BR (%)
	BD-BR (%)	BD-PSNR (dB)	BD-BR (%)	BD-PSNR (dB)	BD-BR (%)	BD-PSNR (dB)	BD-BR (%)	BD-PSNR (dB)	
longdress	-79.65	10.50	-62.23	9.41	-31.12	0.60	-21.50	0.41	-92.37
loot	-83.68	10.11	-73.83	8.89	-28.36	0.78	7.66	-0.07	-43.36
queen	-87.64	11.64	-78.17	9.99	-17.12	0.49	-5.38	0.26	-72.59
redandblack	-79.35	9.35	-54.33	8.28	-32.16	0.87	-21.56	0.21	-50.58
soldier	-80.12	10.09	-68.73	8.86	-41.38	1.05	-29.99	0.63	-56.24
Average	-82.09	10.34	-67.46	9.08	-30.03	0.76	-14.15	0.29	-63.03

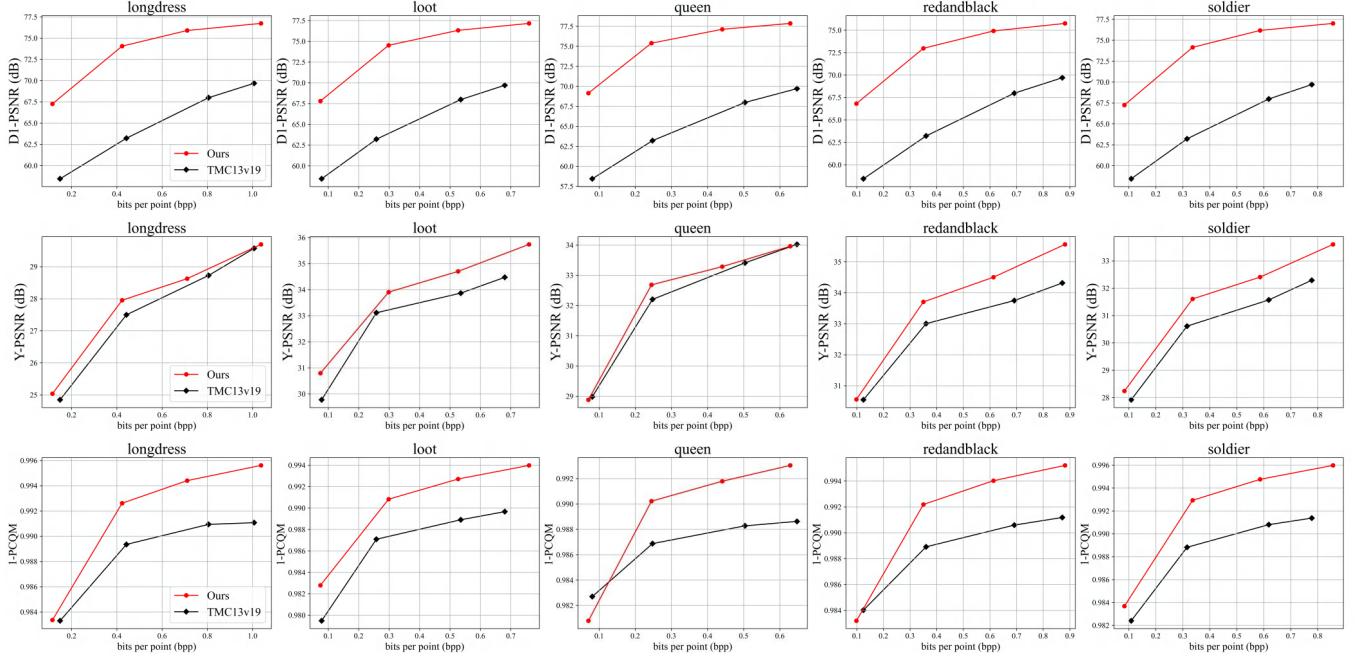


Fig. 11: **R-D Curves of DeepPCC and G-PCC (TMC13v19) for Unified Lossy Solution.** Distortion measures using D1 PSNR, Y PSNR, and PCQM are provided for comparison. The horizontal axis represents the total bit rates. It is worth noting that we ensure DeepPCC and G-PCC have comparable geometric bit rates and comparable attribute bit rates for a fair comparison.

methods but just serve as a reference to provide an intuitive understanding of their performance. On the other hand, since both G-PCC and DeepPCC currently have no accurate rate control technique, bitrate mismatch is inevitable.

We evaluate the restored geometry quality using D1 and D2 metrics and the restored attribute quality using Y and YUV PSNR. Furthermore, we employ the PCQM metric to measure the reconstructed quality of both geometry and attribute. As shown in Table VI, DeepPCC outperforms the G-PCC TMC13v19 by 82.09% (67.46%) BD-BR in D1 (D2) measurement, 30.03% (14.15%) BD-BR in Y (YUV) space, respectively. Regarding the overall PCQM that measures the superimposed distortion from both geometry and attributes, DeepPCC attains 63.03% BD-BR gains over the TMC13v19. Corresponding R-D curves are provided in Fig. 11.

The reconstructed point clouds from different methods are visualized in Fig. 1 for comparison. We observe that the point cloud reconstructed by G-PCC contains abundant compression noise, blocking artifacts, and black holes, which severely deteriorate the visual quality. The blocking artifacts are mainly triggered by the RAHT and quantization algorithms utilized in the G-PCC attribute compression process, while

the black holes come from the absence of geometry coordinates after G-PCC geometry compression. In contrast, the proposed DeepPCC provides much better visual representation with smooth surfaces and clear edges, yielding more visually pleasing reconstructions. We observe that some regions in our reconstruction look smooth. This is because we employ the MSE loss to train DeepPCAC. This problem can be largely resolved through the use of a different loss such as MPED [58], which will be explored in our future work.

We further conduct experiments to evaluate the performance of DeepPCC on the Microsoft Voxelized Upper Bodies (MVUB) [59] dataset, which consists of five point clouds: *Andrew*, *David*, *Phil*, *Ricardo*, and *Sarah*. Corresponding R-D curves are presented in Fig. 12. It is shown that our DeepPCGC outperforms G-PCC by a large margin, providing a superior foundation for the subsequent PCAC processing. Benefiting from the well-reconstructed geometry, our DeepPCAC significantly outperforms G-PCC. As a result, the overall PCQM of DeepPCC, which jointly assesses the quality of geometry and attribute, surpasses that of the G-PCC remarkably.

The comparative study to the JPEG VM is supplemented in a companion material due to space limitation.

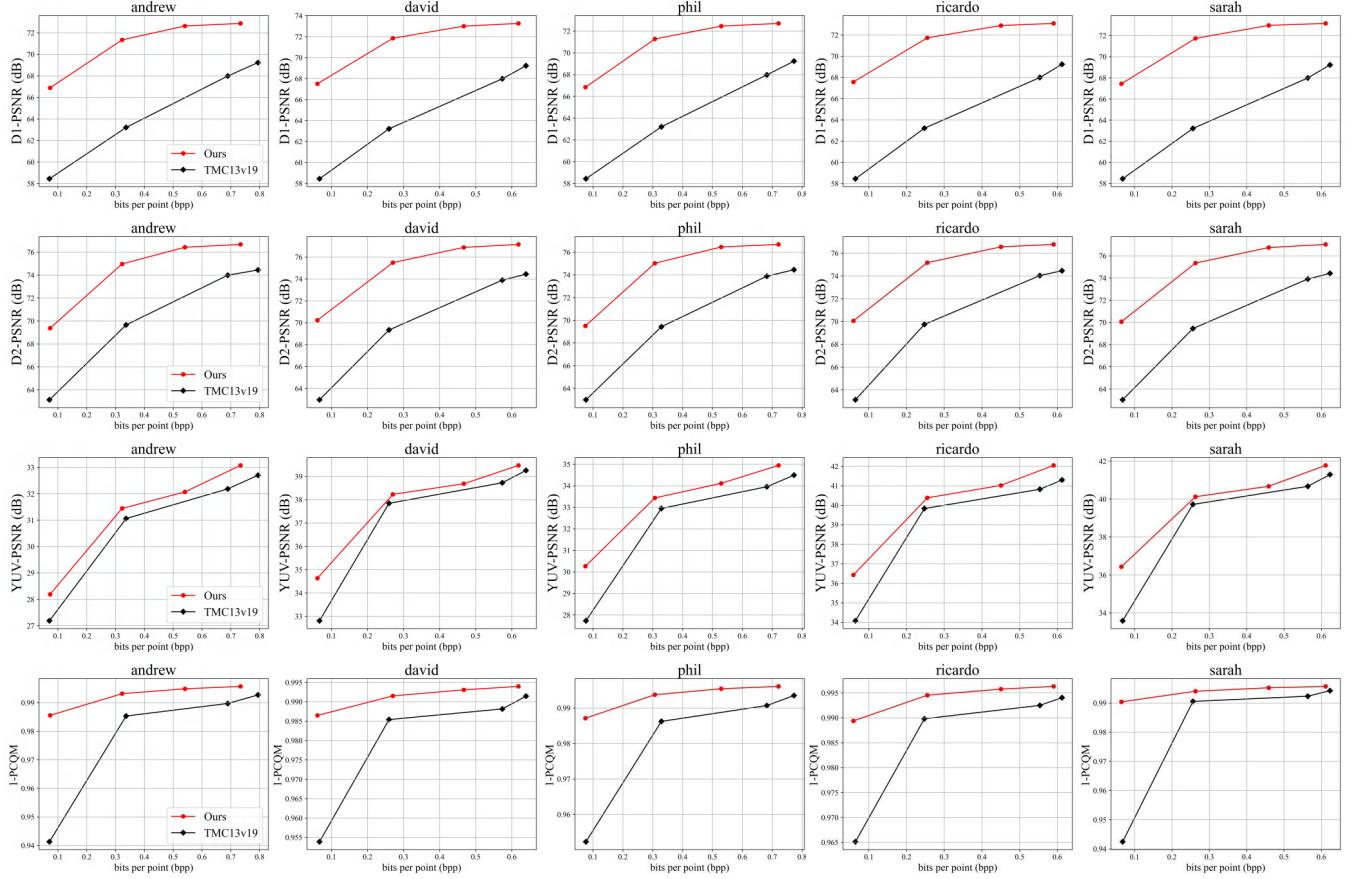


Fig. 12: **R-D Comparison using MVUB Dataset.** Our DeepPCC models are trained using ShapeNet. We ensure DeepPCC and G-PCC have comparable geometric bit rates and comparable attribute bit rates for a fair comparison.

E. Space and Time Complexity

Our DeepPCGC model uses 1.26 M parameters and its model size is 5.30 MB. Our DeepPCAC model has 25.63 M parameters, of which the Attribute Filter uses 1.51 M parameters. Its model size is 119.0 MB, of which the Attribute Filter costs 6 MB.

The runtime of DeepPCC is presented in Table VII. It is observed that its decoding time significantly exceeds its encoding time. This is attributed to the sequential decoding process involved in our entropy coding using the autoregressive model. It is worth noting that using the entropy model with only hyperprior can reduce the runtime by approximately 90% but at the cost of a significant decline in performance, as we have demonstrated in our previous work [25]. Recently, a multistage context modeling method developed in [60] not only reports competitive performance to the context model using autoregressive neighbors but also demonstrates 6 \times to 8 \times decoding speedup, which can be borrowed to accelerate the DeepPCC entropy decoding in the future.

V. CONCLUDING REMARKS

In this work, a learning-based lossy point cloud compression approach – DeepPCC was proposed and evaluated. Leveraging the powerful representation capability of deep learning techniques, DeepPCC employs a simple autoencoder architecture

TABLE VII: DeepPCC Runtime Summary

Point Cloud	DeepPCGC		DeepPCAC		DeepPCC	
	Enc (s)	Dec (s)	Enc (s)	Dec (s)	Enc (s)	Dec (s)
longdress	2.15	16.68	24	216	23.15	220.68
loot	2.01	15.26	20	182	21.01	188.26
queen	2.49	15.66	22	195	29.4	223.44
redandblack	2.01	14.97	19	175	24.6	192.68
soldier	2.85	21.37	30	306	37.22	333.04
Average	2.30	16.79	23.00	214.80	27.08	231.63

consisting of stacked neural network layers to compactly represent the geometry and attribute information of large-scale point clouds. To best exploit inter-point correlations for efficient compression, we develop a multiscale neighborhood information aggregation method that progressively downscale the input point cloud to generate multi-resolution scales, and at each resolution scale, the combination of sparse convolution and local self-attention is applied to effectively characterize and embed correlations across spatial neighbors. Extensive simulations demonstrate the superior performance of DeepPCC over existing methods, such as MPEG G-PCC and JPEG VM, in compressing individual components (*e.g.*, geometry or attribute) or color point cloud with both geometry and attribute information.

For future work, balancing the bitrates of geometry and attribute components for global optimal is an interesting

- [44] J. Wang, D. Ding, and Z. Ma, “Lossless point cloud attribute compression using cross-scale, cross-group, and cross-color prediction,” in *Data Compression Conference (DCC)*. IEEE, 2023.
- [45] G. Fang, Q. Hu, H. Wang, Y. Xu, and Y. Guo, “3DAC: Learning attribute compression for point clouds,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 14 819–14 828.
- [46] D. Ding, J. Zhang, J. Wang, and Z. Ma, “CARNet: Compression artifact reduction for point cloud attribute,” *arXiv preprint arXiv:2209.08276*, 2022.
- [47] D. Minnen, J. Ballé, and G. D. Toderici, “Joint autoregressive and hierarchical priors for learned image compression,” *Advances in neural information processing systems*, vol. 31, 2018.
- [48] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston, “Variational image compression with a scale hyperprior,” *arXiv preprint arXiv:1802.01436*, 2018.
- [49] T. Chen, H. Liu, Z. Ma, Q. Shen, X. Cao, and Y. Wang, “End-to-end learnt image compression via non-local attention optimization and improved context modeling,” *IEEE Transactions on Image Processing*, vol. 30, pp. 3179–3191, 2021.
- [50] L. Kong, D. Ding, F. Liu, D. Mukherjee, U. Joshi, and Y. Chen, “Guided CNN restoration with explicitly signaled linear combination,” in *2020 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2020, pp. 3379–3383.
- [51] 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.
- [52] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft COCO: Common objects in context,” in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [53] Z. Alexandre, D. Graziosi, and T. Ali, “Dataset split for AI-PCC,” *ISO/IEC JTC 1/SC 29/WG 7 m58754*, Jan. 2022.
- [54] WG7, MPEG 3D Graphics Coding, “Description of exploration experiment 5.3 on AI-based dynamic pc coding,” *ISO/IEC JTC 1/SC 29/WG 7 N00386*, July 2022.
- [55] E. d’Eon, B. Harrison, T. Myers, and P. A. Chou, “8i voxelized full bodies - a voxelized point cloud dataset,” *ISO/IEC JTC1/SC29 Joint WG11/WG1 (MPEG/JPEG) m38673/M72012*, May 2016.
- [56] X. Yi, L. Yao, and W. Ziyu, “Owlpii dynamic human mesh sequence dataset,” *ISO/IEC JTC1/SC29/WG11 (MPEG/JPEG) m41658*, October, 2017.
- [57] C. Choy, J. Gwak, and S. Savarese, “4D spatio-temporal ConvNets: Minkowski convolutional neural networks,” *2019 IEEE/CVF CVPR*, pp. 3070–3079, 2019.
- [58] Q. Yang, Y. Zhang, S. Chen, Y. Xu, J. Sun, and Z. Ma, “MPED: Quantifying point cloud distortion based on multiscale potential energy discrepancy,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [59] C. Loop, Q. Cai, S. O. Escalano, and P. A. Chou, “Microsoft voxelized upper bodies-a voxelized point cloud dataset,” *ISO/IEC JTC1/SC29 Joint WG11/WG1 (MPEG/JPEG) input document m38673 M*, vol. 72012, p. 2016, 2016.
- [60] M. Lu and Z. Ma, “High-efficiency lossy image coding through adaptive neighborhood information aggregation,” *arXiv preprint arXiv:2204.11448*, 2022.
- [61] W. Zhu, Y. Xu, D. Ding, Z. Ma, and M. Nilsson, “Lossy point cloud geometry compression via region-wise processing,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 12, pp. 4575–4589, 2021.
- [62] S. Perry, “JPEG pleno point cloud coding common test conditions v3.2,” in *87th JPEG Meeting, WG1N87037*, 2020.
- [63] A. F. Guarda, N. M. Rodrigues, M. Ruivo, L. Coelho, A. Seleem, and F. Pereira, “IT/IST/IPLeiria response to the call for proposals on jpeg pleno point cloud coding,” *arXiv preprint arXiv:2208.02716*, 2022.

SUPPLEMENTARY MATERIAL

A. Detailed Network Architecture of DeepPCGC and DeepPCAC

The network layers of DeepPCC for geometry and attribute compression are detailed in Fig. 1a and Fig. 1b, respectively.

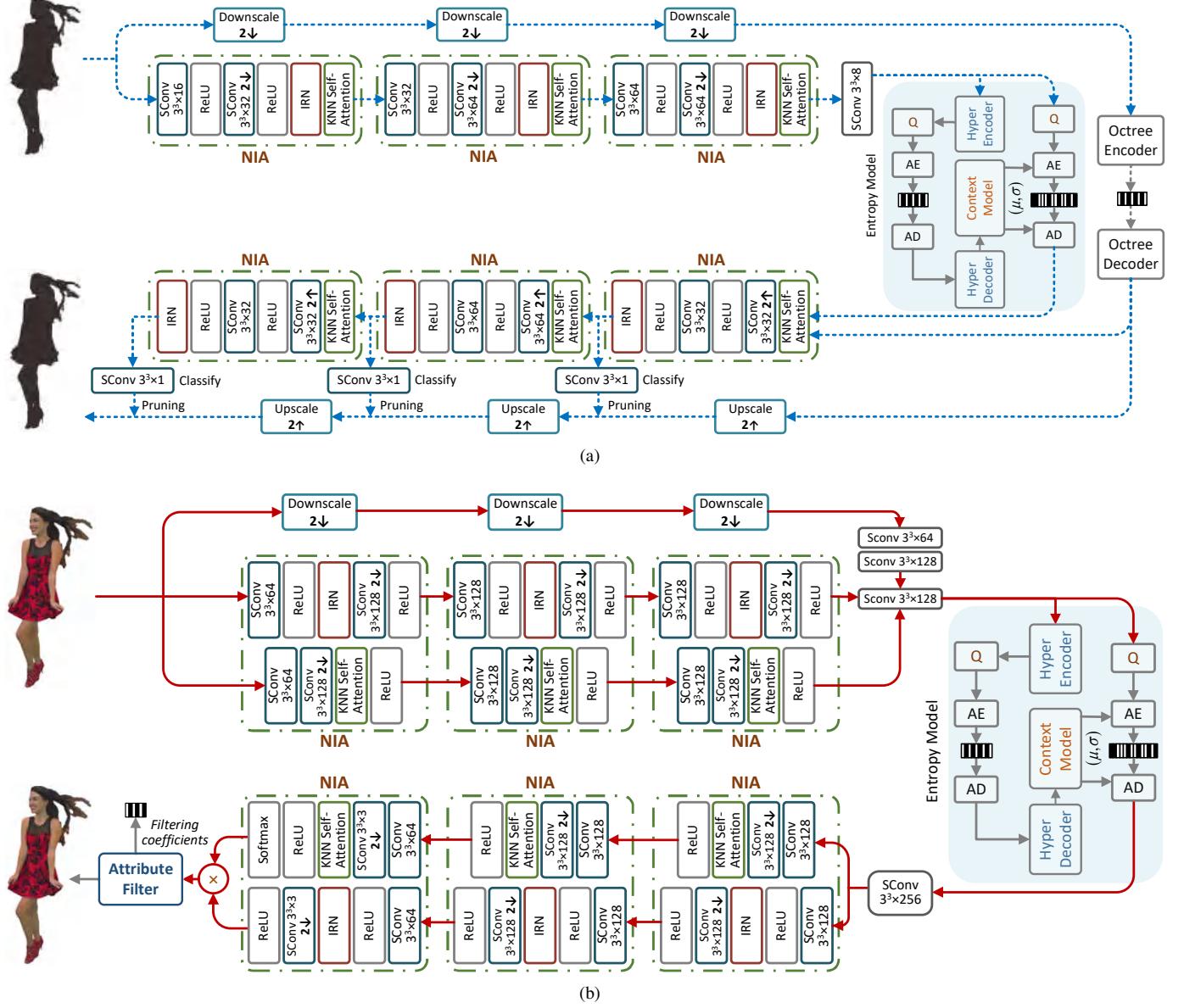


Fig. 1: Network Architecture for Point Cloud Geometry Compression (PCGC) and Point Cloud Attribute Compression (PCAC) in DeepPCC. (a) DeepPCGC. (b) DeepPCAC.

B. NIA Module

Previous discussions have demonstrated that the incorporation of NIA modules significantly contributes to the performance of DeepPCC. In this section, we will provide more details on the NIA which consists of IRN and KNN Self-Attention.

IRN. The IRN module is essentially an inception-residual network. As depicted in Fig. 3a of our main manuscript, besides the skip connection branch, the IRN has another two parallel branches: one stacks three SConvs with the kernel size of 1×1 , $3 \times 3 \times 3$, and $1 \times 1 \times 1$, and the other stacks two $3 \times 3 \times 3$ SConvs. As such, the IRN is capable of capturing both long-range and short-range reception fields, facilitating multiple-level information collection. The skip connection additionally helps the network to converge.

Local Self-Attention. In DeepPCC, the quantized latent feature vectors obtained at the bottleneck, which contain the scale-wise aggregated information, are critical to compression performance as they are used to hierarchically reconstruct a point cloud at the decoder. Previously approaches, such as stacked sparse convolutions were mainly used for feature analysis and

aggregation. However, their information aggregation efficiency was often limited because preset fixed convolutions are unable to effectively capture sufficient points, particularly in geometry regions with low point density, and thus fail to characterize the dynamic content (*e.g.*, regions with unevenly distributed points) [61].

The DeepPCC hence introduces the local Self-Attention mechanism to tackle the above issues. It basically comprises two key components: 1) dynamic local neighborhood construction for each point using KNN search, and 2) adaptive information aggregation using Transformer-based Self-Attention over the K nearest neighbors.

Local Neighborhood Construction Using KNN Search. We introduce a KNN module to instantaneously collect K nearest neighbors of each point for local neighborhood generation. As seen, this local neighborhood is dynamically conditioned on each individual point.

As presented in Fig. 3b in our main manuscript, at the beginning of each Transformer module, K neighboring points of the current point along with their coordinates and features are gathered and combined as F_X'' of $N \times K \times d$ for further processing. Here, d denotes the number of output feature maps. Compared to previous works which input the entire point cloud into networks for processing, our KNN-based method greatly reduces the memory requirement and computational complexity. On the other hand, the nearest neighbors generally exhibit stronger correlations with the current point, whereas points located at a greater distance may present less correlation or even negative contributions. As a result, the selection of K will impact the coding performance and complexity. In this study, we set $k = 16$.

Local Information Aggregation Using Self-Attention. Subsequently, latent features including F_X and F_X'' are both sent into linear layers for computing the Q , K , and V as the conventional Transformer does, shown in Fig. 3c in our main manuscript. For every point, an attention map will be calculated. By applying the attention map to K local neighbor features, we attain the attention-weighted feature at a size of $1 \times d$ of each point. As such, the information from valid K nearest neighbors is aggregated and embedded compactly as latent features of a point. The feature vector of each point describes its spatial neighborhood embedding. For all points of the input PCG, we finally obtain attentive features with the size of $N \times d$. As illustrated in Fig. 1, features are progressively aggregated from one scale to another and eventually encoded at the bottleneck.

C. Experimental Settings of Fig. 1

Our DeepPCC is a unified lossy solution including DeepPCGC for geometry and DeepPCAC for attribute compression, which is trained using the ShapeNet dataset. The models of JPEG VM are from its opensource website. The latest G-PCC version (TMC13v19) including G-PCC (octree) for geometry compression and G-PCC (RAHT) for attribute compression is used for comparison. To ensure a fair comparison, we keep comparable bit rates between the DeepPCGC and G-PCC (octree) and comparable bit rates between the DeepPCAC and G-PCC (RAHT). Since JPEG VM devises an end-to-end model for monolithic geometry and attribute compression, the respective bit rates for geometry and attribute are unavailable. Specifically, for Fig. 1 in our main manuscript, we detail the bit rate allocation results of G-PCC and DeepPCC in Table I.

TABLE I: Bit Rate Allocation for Geometry and Attribute in G-PCC and DeepPCC

Point Cloud	Method	Geometry bit rate (bpp)	Attribute bit rate (bpp)	Total bit rate (bpp)
redandblack	G-PCC TMC13v6	0.5482	0.0921	0.6403
	G-PCC TMC13v19	0.2518	0.1084	0.3602
	DeepPCC	0.2300	0.1201	0.3501
solider	G-PCC TMC13v6	0.5464	0.0637	0.6103
	G-PCC TMC13v19	0.2311	0.0846	0.3157
	DeepPCC	0.2320	0.1043	0.3363

D. Comparison with JPEG VM

We also conduct a comparison between our proposed DeepPCC and the JPEG VM. We additionally select four point clouds sequences from JPEG Pleno Point Cloud Test database [62], including RWT136, RWT395, RWT503b, and RWT529 (as visualized in Fig. 2).

This comparison is only intended as a reference since the JPEG VM is still under development. We limit our comparison to the JPEG test sequences since the 8iVFB sequences are used as the training dataset of JPEG VM. The comparison results, as shown in Fig. 3, reveal that our DeepPCC outperforms the JPEG VM in various distortion measurements, such as D1, YUV PSNR, and PCQM.

E. Open Source

DeepPCC is built upon several pioneering explorations as listed in Table II, which will be made publicly accessible soon at <https://github.com/3dpcc>.

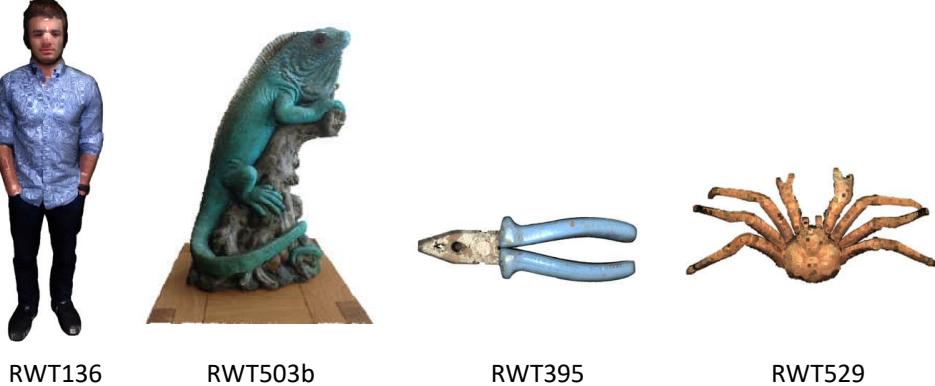


Fig. 2: Test Sequences from JPEG Pleno Point Cloud Test database. Four point clouds are randomly selected for testing.

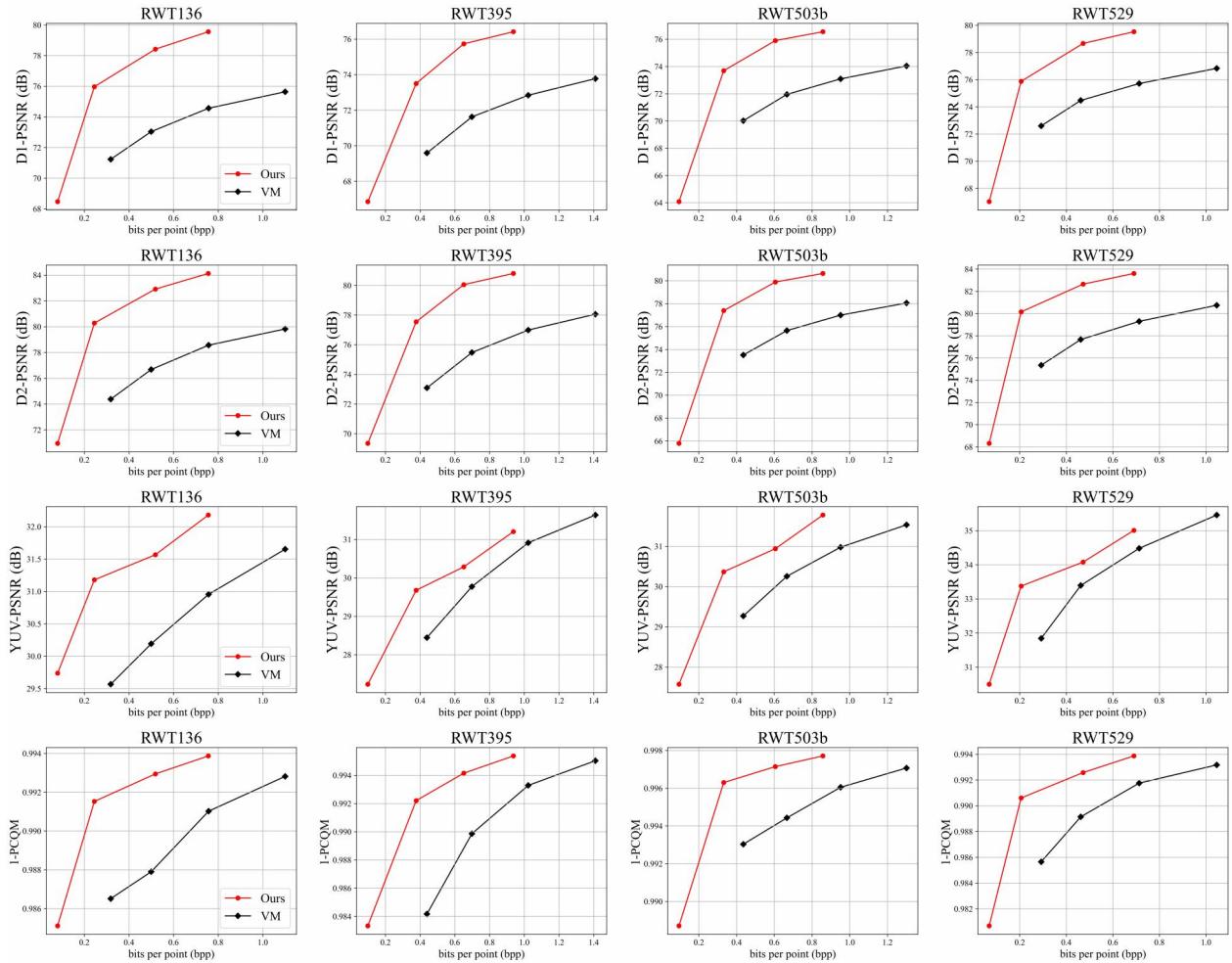


Fig. 3: R-D Curve Comparison with JPEG VM in JPEG Test Sequences. The JPEG VM models are from the publicly accessible website and they are pretrained following JPEG training conditions [62]. Our DeepPCC models are trained using ShapeNet.

TABLE II: Open-source Examples of Point Cloud Compression Implementations

Codec	Component	Description	Open Source
PCGCv1 [23]	Geometry	End-to-End Learning	https://njuvision.github.io/PCGCv1/
PCGCv2 [25]	Geometry	End-to-End Learning	https://njuvision.github.io/PCGCv2/
SparsePCGC [22]	Geometry	End-to-End Learning	https://github.com/NJUVISION/SparsePCGC
JPEG Pleno VM [63]	Geometry and Attribute	End-to-End Learning	https://github.com/aguarda/JPEG-VM
G-PCC	Geometry and Attribute	Classical Rules-based	https://github.com/MPEGGroup/mpeg-pcc-tmc13
DeepPCC	Geometry and Attribute	End-to-End Learning	https://github.com/3dpcc