
Diff-PCC: Diffusion-based Neural Compression for 3D Point Clouds

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Stable diffusion networks have emerged as a groundbreaking development for
2 their ability to produce realistic and detailed visual content. This characteristic
3 renders them ideal decoders, capable of producing high-quality and aesthetically
4 pleasing reconstructions. In this paper, we introduce the first diffusion-based point
5 cloud compression method, dubbed Diff-PCC, to leverage the expressive power of
6 the diffusion model for generative and aesthetically superior decoding. Different
7 from the conventional autoencoder fashion, a dual-space latent representation
8 is devised in this paper, in which a compressor composed of two independent
9 encoding backbones is considered to extract expressive shape latents from distinct
10 latent spaces. At the decoding side, a diffusion-based generator is devised to
11 produce high-quality reconstructions by considering the shape latents as guidance
12 to stochastically denoise the noisy point clouds. Experiments demonstrate that the
13 proposed Diff-PCC achieves state-of-the-art compression performance (e.g., 7.711
14 dB BD-PSNR gains against the latest G-PCC standard at ultra-low bitrate) while
15 attaining superior subjective quality. Source code will be made publicly available.

16 1 Introduction

17 Point clouds, composed of numerous discrete points with coordinates (x, y, z) and optional attributes,
18 offer a flexible representation of diverse 3D shapes and are extensively applied in various fields such
19 as autonomous driving [8], game rendering [35], robotics [7], and others. With the rapid advancement
20 of point cloud acquisition technologies and 3D applications, effective point cloud compression
21 techniques have become indispensable to reduce transmission and storage costs.

22 1.1 Background

23 Prior to the widespread adoption of deep learning techniques, the most prominent traditional point
24 cloud compression methods were the G-PCC [39] and V-PCC [40] proposed by the Moving Picture
25 Experts Group(MPEG). G-PCC compresses point clouds by converting them into a compact tree
26 structure, whereas V-PCC projects point clouds onto a 2D plane for compression. In recent years,
27 numerous deep learning-based methods have been proposed [50, 45, 11, 12, 7, 30, 46, 14, 42],
28 which primarily employ the Variational Autoencoder (VAE) [1, 2] architecture. By learning a prior
29 distribution of the data, the VAE projects the original input into a higher-dimensional latent space,
30 and reconstructs the latent representation effectively using a posterior distribution. However, previous
31 VAE-based point cloud compression architectures still face recognized limitations: 1) Assuming a
32 single Gaussian distribution $N(\mu, \sigma^2)$ in the latent space may prove inadequate to capture the intricate
33 diversity of point cloud shapes, yielding blurry and detail-deficient reconstructions [56, 10]; 2) The
34 Multilayer Perceptron (MLP) based decoders [50, 45, 11, 12, 46] suffer from feature homogenization,
35 which leads to point clustering and detail degradations in the decoded point cloud surfaces, lacking the

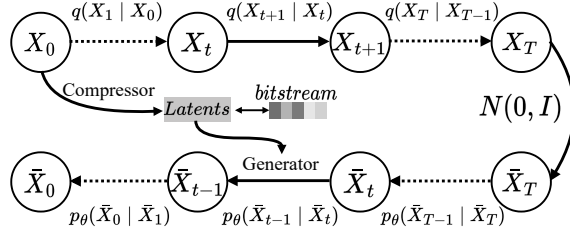


Figure 1: Diff-PCC pipeline. X_t and \bar{X}_t represents the t th original point cloud and noisy point cloud, respectively; p refers to the forward process and q refers to the reverse process; $N(0, I)$ means the pure noise. Entropy model and arithmetic coding is omitted for a concise explanation.

36 ability to produce high-quality reconstructions. Recently, Diffusion models (DMs) [5] have attracted
 37 considerable attention in the field of generative modeling [34, 48, 41, 19] due to their outstanding
 38 performance in generating high-quality samples and adapting to intricate data distributions, thus
 39 presenting a novel and exciting opportunity within the domain of neural compression [33, 44, 25].
 40 By generating a more refined and realistic 3D point cloud shape, DMs offer a distinctive approach to
 41 reduce the heavy dependence of reconstruction quality on the information loss of bottleneck layers.

42 1.2 Our Approach

43 Building on the preceding discussion, we introduce Diff-PCC, a novel lossy point cloud compression
 44 framework that leverages diffusion models to achieve superior rate-distortion performance with
 45 exceptional reconstruction quality. Specifically, to enhance the representation ability of simplistic
 46 Gaussian priors in VAEs, this paper devises a dual-space latent representation that employs two
 47 independent encoding backbones to extract complementary shape latents from distinct latent spaces.
 48 At the decoding side, a diffusion-based generator is devised to produce high-quality reconstructions by
 49 considering the shape latents as guidance to stochastically denoise the noisy point clouds. Experiments
 50 demonstrate that the proposed Diff-PCC achieves state-of-the-art compression performance (e.g.,
 51 7.711 dB BD-PSNR gains against the latest G-PCC standard at ultra-low bitrate) while attaining
 52 superior subjective quality.

53 1.3 Contribution

54 Main contributions of this paper are summarized as follows:

- 55 • We propose Diff-PCC, a novel diffusion-based lossy point cloud compression framework.
 56 To the best of our knowledge, this study presents *the first* exploration of diffusion-based
 57 neural compression for 3D point clouds.
- 58 • We introduce a dual-space latent representation to enhance the representation ability of the
 59 conventional Gaussian priors in VAEs, enabling the Diff-PCC to extract expressive shape
 60 latents and facilitate the following diffusion-based decoding process.
- 61 • We devise an effective diffusion-based generator to produce high-quality noises by consider-
 62 ing the shape latents as guidance to stochastically denoise the noisy point clouds.

63 2 Related Work

64 2.1 Point Cloud Compression

65 Classic point cloud compression standards, such as G-PCC, employ octree[29] to compress point
 66 cloud geometric information. In recent years, inspired by deep learning methods in point cloud
 67 analysis[26, 27] and image compression[1, 2, 22], researchers have turned their attention to learning-
 68 based point cloud compression. Currently, point cloud compression methods can be primarily divided
 69 into two branches: voxel-based and point-based approaches. Voxel-based methods further branch into

70 sparse convolution[36, 37, 38, 49, 51, 52] and octree[9, 24, 31]. Among them, sparse convolution de-
 71 rives from 2D-pixel representations but optimizes for voxel sparsity. On the other hand, octree-based
 72 methods, utilize tree structures to eliminate redundant voxels, representing only the occupied ones.
 73 Point-based methods[11, 50, 45, 46] are draw inspiration from PointNet [26], utilizing symmetric
 74 operators (max pooling, average pooling, attention pooling) to handle permutation-invariant point
 75 clouds and capture geometric shapes. For compression, different quantization operations categorize
 76 point cloud compression into lossy and lossless types. In this paper, we focus on lossy compression
 77 to achieve higher compression ratios by sacrificing some precision in the original data.

78 2.2 Diffusion Models for Point Cloud

79 Recently, diffusion models have ignited the image generation field[58, 17, 32], inspiring researchers
 80 to explore their potential in point cloud applications. DPM[20] pioneered the introduction of diffusion
 81 models in this domain. Starting from DPM, PVD[57] combines the strengths of point cloud and
 82 voxel representations, establishing a baseline based on PVCNN. LION[47] employs two diffusion
 83 models to separately learn shape representations in latent space and point representations in 3D
 84 space. Dit-3D[23] innovates by integrating transformers into DDPM, directly operating on voxelized
 85 point clouds during the denoising process. PDR[21] employs diffusion model twice during the
 86 process of generating coarse point clouds and refined point clouds. Point-E[] utilizes three diffusion
 87 models for the following processes: text-to-image generation, image-to-point cloud generation, and
 88 point cloud upsampling. PointInfinity[13] utilizes cross-attention mechanism to decouple fixed-size
 89 shape latent and variable-size position latent, enabling the model to train on low-resolution point
 90 clouds while generating high-resolution point clouds during inference. DiffComplete[4] enhances
 91 control over the denoising process by incorporating ControlNet[53], achieving new state-of-the-art
 92 performances. These advancements demonstrate the promise of DMs in point cloud generation tasks,
 93 which motivates our exploring its applicability in point cloud compression. Our research objective is
 94 to explore the effective utilization of diffusion models for point cloud compression while preserving
 95 its critical structural features.

96 3 Method

97 Figure 1 illustrates the pipeline of the proposed Diff-PCC, which can also represent the general work-
 98 flow of diffusion-based neural compression. A concise review for Denoising Diffusion Probabilistic
 99 Models (DDPMs) and Neural Network (NN) based point cloud compression is first provided in
 100 Sec. 3.1; The proposed Diff-PCC is detailed in Sec. 3.2.

101 3.1 Preliminaries

102 Denoising Diffusion Probabilistic Models (DDPMs) comprise two Markov chains of length T :
 103 diffusion process and denoising process. Diffusion process adds noise to clean data \mathbf{x}_0 , resulting in
 104 a series of noisy samples $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$. When T is large enough, $\mathbf{x}_T \sim N(0, \mathbf{I})$. The denoising
 105 process is the reverse process, gradually removing the noise added during the diffusion process. We
 106 formulate them as follows:

$$q(\mathbf{x}_1, \dots, \mathbf{x}_T | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1}), \text{ where } q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}) \quad (1)$$

$$p_{\theta}(\mathbf{x}_0, \dots, \mathbf{x}_{T-1} | \mathbf{x}_T) = \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t), \text{ where } p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I}) \quad (2)$$

107 where β is a hyperparameter representing noise level. $t \sim \text{Unif}\{1, \dots, T\}$ represents time step. Via
 108 reparameterization trick, we can sample from $q(\mathbf{x}_t | \mathbf{x}_{t-1})$ and $p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$ as following:

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \boldsymbol{\epsilon} \quad (3)$$

$$\mathbf{x}_{t-1} = \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t) + \sigma_t \boldsymbol{\epsilon} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}} \beta_t \boldsymbol{\epsilon} \quad (4)$$

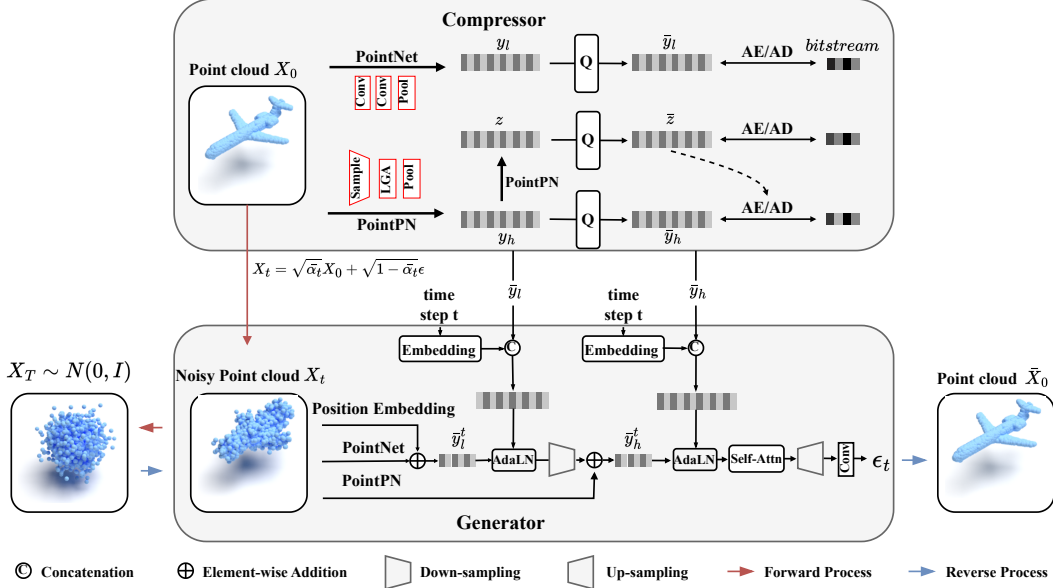


Figure 2: Detailed Structure of the Utilized Compressor and Generator. y_l and y_h refer to the low-frequency shape latent and high-frequency detail latent, respectively; z means hyperprior latent; Q refers to the quantization; AE and AD represents the arithmetic encoding and decoding.

109 where $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$, ϵ denotes random noise sampled from $\mathcal{N}(0, \mathbf{I})$. Note that
 110 $\epsilon_\theta(\mathbf{x}_t, t)$ is a neural network used to predict noise during the denoising process, and \mathbf{x}_t can be
 111 directly sampled via $x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon$.

112 DDPMs train the reverse process by optimizing the model parameters θ through noise distortion. The
 113 loss function $L(\theta, \mathbf{x}_0)$ is defined as the expected squared difference between the predicted noise and
 114 the actual noise, with the mathematical expression as follows:

$$L(\theta, \mathbf{x}_0) = \mathbf{E}_{t, \epsilon} \|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2 \quad (5)$$

115 3.2 DIFF-PCC

116 3.2.1 Overview

117 As shown in Fig. 2, two key components, i.e., compressor and generator, are respectively utilized
 118 in the diffusion process and denoising process. In Diff-PCC, the diffusion process is identified as
 119 the encoding, in which a compressor extracts latents from the point cloud and compresses latents
 120 into bitstreams; at the decoding side, the generator accepts the latents as a condition and gradually
 121 restoring point cloud shape from noisy samples.

122 3.2.2 Dual-Space Latent Encoding

123 Several research have demonstrated that a simplistic Gaussian distribution in the latent space may
 124 prove inadequate to capture the complex visual signals [56, 3, 6, 10]. Although previous works have
 125 proposed to solve these problems using different technologies such as non-gaussian prior [15] or
 126 coupling between the prior and the data distribution [10], these techniques may not be able to directly
 127 employed on neural compression tasks.

128 In this paper, a simple yet effective compressor is introduced, which composed of two independent
 129 encoding backbones to extract expressive shape latents from distinct latent spaces. Motivated by
 130 PointPN [55], which excels in capturing high-frequency 3D point cloud structures characterized by
 131 sharp variations, we design a dual-space latent encoding approach that utilizes PointNet to extract
 132 low-frequency shape latent and leverages PointPN to characterize complementary latent from high
 133 frequency domain. Let x be the original input point cloud, we formulate the above process as:

$$\{y_l, y_h\} = \{E_l(x), E_h(x)\} \quad (6)$$

134 where $y_l \in \mathbb{R}^{1 \times C}$ and $y_h \in \mathbb{R}^{S \times C}$ represent the low-frequency and high-frequency latent features,
 135 respectively; E_l and E_h refer to the PointNet and PointPN backbones, respectively. Next, the
 136 quantization process Q is applied on the obtained features \bar{y}_l and \bar{y}_h , i.e.,

$$\{\bar{y}_l, \bar{y}_h\} = \{Q(y_l), Q(y_h)\} \quad (7)$$

137 where function Q refers to the operation of adding uniform noise during training [1] and the rounding
 138 operation during test.

139 Then, fully factorized density model [1] and the hyperprior density model [2] are employed to fit the
 140 distribution of quantized features \bar{y}_l and \bar{y}_h , respectively. Particularly, the hyperprior density model
 141 $p_\varphi(\bar{y}_h)$ can be described as:

$$p_\varphi(\bar{y}_h) = \left(N(\mu, \sigma^2) * \mathcal{U}\left(-\frac{1}{2}, \frac{1}{2}\right) \right) (\bar{y}_h) \quad (8)$$

142 where $\mathcal{U}\left(-\frac{1}{2}, \frac{1}{2}\right)$ refers to the uniform noise ranging from $-\frac{1}{2}$ to $\frac{1}{2}$; $N(\mu, \sigma^2)$ refers to the normal
 143 distribution with expectation μ and standard deviation σ , which can be further estimated by a
 144 hyperprior encoder E_{hyper} and decoder D_{hyper} :

$$(\mu, \sigma^2) = D_{hyper}(\bar{z}) = D_{hyper}(Q(z)) = D_{hyper}(Q(E_{hyper}(y_h))) \quad (9)$$

145 In this way, a triplet containing quantized low-frequency feature \bar{y}_l , quantized high-frequency feature
 146 \bar{y}_h , and quantized hyperprior \bar{z} will be compressed into three separate streams. Let $p(\cdot)$ and $p_{(\dots)}(\cdot)$
 147 respectively represents the actual distribution and estimated distribution of latent features, then the
 148 bitrate \mathcal{R} can be estimated as follows:

$$\mathcal{R} = \mathbb{E}_{\bar{y}_l \sim p(\bar{y}_l)} [-\log_2 p_\theta(\bar{y}_l)] + \mathbb{E}_{\bar{y}_h \sim p(\bar{y}_h)} [-\log_2 p_\varphi(\bar{y}_h)] + \mathbb{E}_{\bar{z} \sim p(\bar{z})} [-\log_2 p_\phi(\bar{z})] \quad (10)$$

149 3.2.3 Diffusion-based Generator

150 The generator takes noisy point cloud x_t at time t and necessary conditional information C as input.
 151 We hope generator to learn positional distribution F of x_t and fully integrate F with C to predict
 152 noise ϵ_t at time t . In this paper, we consider all information that could potentially guide the generator
 153 as conditional information, including time t , class label l , noise coefficient β_t , and decoded latent
 154 features (\bar{y}_l and \bar{y}_h).

155 DiffComplete [4] uses ControlNet [54] to achieve refined noise generation. However, the denoiser of
 156 DiffComplete is a 3D-Unet, adapted from its 2D version [16]. This structure is not suitable for our
 157 method, because we directly deal with points, instead of voxels. We embraced this idea and specially
 158 designed a hierarchical feature fusion mechanism to adapt to our method. Note that 3D-Unet can
 159 directly downsample features F through 3D convolution with a stride greater than one. It is very
 160 complex for point-based methods to achieve equivalent processing. Therefore, we did not replicate
 161 the same structure as DiffComplete does, but directly used AdaLN to inject conditional information,
 162 formulated as:

$$AdaLN(F_{in}, C) = Norm(F_{in}) \odot Linear(C) + Linear(C) \quad (11)$$

163 where F_{in} denotes the original features in the Generator and C denotes the condition information.

164 Now we detail the structure: First, we need to exact the shape latent of noise point cloud x_t and we
 165 choose PointNet for structural consistency. However, in the early stages of the denoising process,
 166 x_t lacks a regular surface shape for the generator to learn. Therefore, we adopt the suggestion from
 167 PDR [23], adding positional encoding to each noise point so that the generator can understand the
 168 absolute position of each point in 3D space. Then we inject shape latent \bar{y}_l from the compressor via
 169 ADA LN. We formulate the above process as:

$$F_{x_t} = PointNet(x_t) + PE(x_t) \quad (12)$$

$$F'_{x_t} = AdaLN(F_{x_t}, C) \quad (13)$$

170 Next, we need to fuse high-frequency features. We extract the local high-frequency features of x_t
 171 using PointPN and add them to F from the previous step, Then we inject the high-frequency features
 172 from the compressor via AdaLN. We use K-Nearest Neighbor (KNN) operation to partition locally

173 and set the number of neighbor points to 8, which allows the generator to learn local details. We
 174 formulate the above process as:

$$F' = PointPN(x_t) + FPS(F_{in}) \quad (14)$$

$$F_{out} = AdaLN(F', C) \quad (15)$$

175 After that, we use the self-attention mechanism to interact with information from different local areas.
 176 And through a feature up-sampling module, we generate features for n points. Finally, we output
 177 noise through a linear layer. We formulate the above process as:

$$F' = SA(F_{in}) \quad (16)$$

$$F'' = UP(F') \quad (17)$$

$$\epsilon_t = Linear(F'') \quad (18)$$

178 3.2.4 Training Objective

179 We follow the conventional rate-distortion trade-off as our loss function as follows:

$$\mathcal{L} = \mathcal{D} + \lambda \mathcal{R} \quad (19)$$

180 where \mathcal{D} refers to the evaluated distortion; \mathcal{R} represents bitrate as shown in Eq. 10; λ serves as the
 181 balance the distortion and bitrate. Specifically, a combined form of distortion \mathcal{D} is used in this paper,
 182 which considers both intermediate noises $(\epsilon, \bar{\epsilon})$ and global shapes (x_0, \bar{x}_0) :

$$\mathcal{D} = \mathcal{D}_{MSE}(\epsilon, \bar{\epsilon}) + \gamma \mathcal{D}_{CD}(x_0, \bar{x}_0) \quad (20)$$

183 where \mathcal{D}_{MSE} denotes the Mean Squared Error (MSE) distance; \mathcal{D}_{CD} refers to the Chamfer Distance;
 184 γ means the weighting factor. Here, the overall point cloud shape is additively supervised under the
 185 Chamfer Distance $\mathcal{D}_{CD}(x_0, \bar{x}_0)$ to provide a global optimization. The following function is utilized
 186 to predict the reconstructed point cloud \bar{x}_0 in practice:

$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(x_t, t, c)) \quad (21)$$

187 where $\bar{\alpha}_t$ means the noise level; x_t refers to the noisy point cloud at time step t ; ϵ_θ denotes the
 188 predicted noise from the generator; c represent the conditional information we inject into the generator.

189 4 Experiments

190 4.1 Experimental Setup

191 **Datasets** Based on previous work, we used ShapeNet as our training set, sourced from [20]. This
 192 dataset contains 51,127 point clouds, across 55 categories, which we allocated in an 8:1:1 ratio for
 193 training, validation, and testing. Each point cloud has 15K points, and following the suggestions from
 194 [28], we randomly select 2K points from each for training. Additionally, we also used ModelNet10
 195 and ModelNet40 as our test sets, sourced from [43]. These datasets contain 10 categories and
 196 40 categories respectively, totaling 10,582 point clouds. During training and testing, we perform
 197 individual normalization on the shape of each point cloud.

198 **Baselines & Metric** We compare our method with the state-of-the-art non-learning-based method:
 199 G-PCC, and the latest learning-based methods from the past two years: IPDAE, PCT-PCC, Following
 200 [45, 46], we use point-to-point PSNR to measure the geometric accuracy and the number of bits per
 201 point to measure the compression ratio.

202 **Implementation** Our model is implemented using PyTorch [27] and CompressAI [4], trained on the
 203 NVIDIA 4090X GPU (24GB Memory) for 80,000 steps with a batch size of 48. We utilize the Adam
 204 optimizer [21] with an initial learning rate of 1e-4 and a decay factor of 0.5 every 30,000 steps, with
 205 β_1 set to 0.9 and β_2 set to 0.999. Since the positional encoding method requires the dimension (dim)
 206 to be a multiple of 6, we designed the bottleneck layer size to be 288. For diffusion, we employ a
 207 cosine preset noise parameter, setting the denoising steps T to 200, which is used for both training
 208 and testing.

Table 1: Objective comparison using BD-PSNR and BD-Rate metrics. G-PCC serves as the anchor. The best and second-best results are highlighted in **bold** and underlined, respectively.

Dataset	Metric	G-PCC	IPDAE	PCT-PCC	Diff-PCC
ShapeNet	BD-Rate (%)	-	-34.594	<u>-87.563</u>	-99.999
	BD-PSNR (dB)	-	+3.518	<u>+8.651</u>	+11.906
ModelNet10	BD-Rate (%)	-	-35.640	-68.899	<u>-56.910</u>
	BD-PSNR (dB)	-	+4.060	+6.333	<u>+5.876</u>
ModelNet40	BD-Rate (%)	-	<u>-53.231</u>	-34.127	-56.451
	BD-PSNR (dB)	-	+4.245	+6.167	<u>+5.350</u>
Avg.	BD-Rate (%)	-	-41.550	<u>-63.530</u>	-71.117
	BD-PSNR (dB)	-	+3.941	<u>+4.384</u>	+7.711
Time (s/frame)	Encoding	0.002	0.004	0.046	0.152
	Decoding	0.001	0.006	0.001	1.913

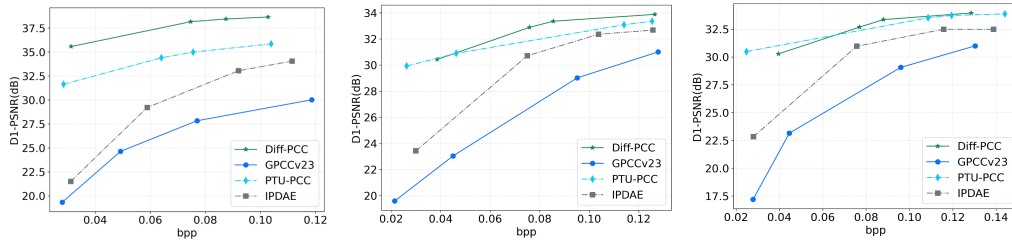


Figure 3: Rate-distortion curves for performance comparison. From left to right: ShapeNet, ModelNet10, and ModelNet40 dataset.

209 4.2 Baseline Comparisons

210 **Objective Quality Comparison** Table 1 shows the quantitative indicators using BD-Rate and BD-
 211 PSNR, and Fig. 3 demonstrates the rate-distortion curves of different methods. It can be seen
 212 that, under identical reconstruction quality conditions, our method achieves superior rate-distortion
 213 performance, conserving between 56% to 99% of the bitstream compared to G-PCC. At the most
 214 minimal bit rates, point of point PSNR of our proposed method surpasses that of G-PCC by 7.711 dB.

215 **Subjective Quality Comparison** Fig 4 presents the ground truth and decoded point clouds from
 216 different methods. We choose three point cloud:airplane, chair, and mug. to be tested across a
 217 comparable bits per pixel (bpp) range. The comparative analysis reveals that at the lowest code rate,
 218 our method preserves the ground truth’s shape information to the greatest extent while simultaneously
 219 achieving the highest Peak Signal-to-Noise Ratio (PSNR).

220 4.3 Ablation Studies

221 We conduct ablation studies to examine the impact of key components in the model. Specifically,
 222 we investigate the effectiveness of low-frequency features, high-frequency features, and the loss
 223 function designed in Sec. 3.2.4. As shown in Table 2, utilizing solely low-frequency features to
 224 guide the reconstruction of the diffusion model results in a 20% reduction in the code rate, along
 225 with a decrease in the reconstruction quality by 0.397dB. This indicates that high-frequency features
 226 play an effective role in guiding the model during the reconstruction process. Conversely, discarding
 227 the low-frequency features, which represent the shape of the point cloud, leads to a reduction in
 228 the code rate and significantly diminishes the reconstruction quality. Therefore, we argue that the
 229 loss of the shape variable is not worth it. Lastly, we ascertain the impact of $\mathcal{D}_{CD}(x_0, \bar{x}_0)$, and the
 230 results indicate that this loss marginally increases the bits per point (bpp) while diminishing the
 231 reconstruction quality.

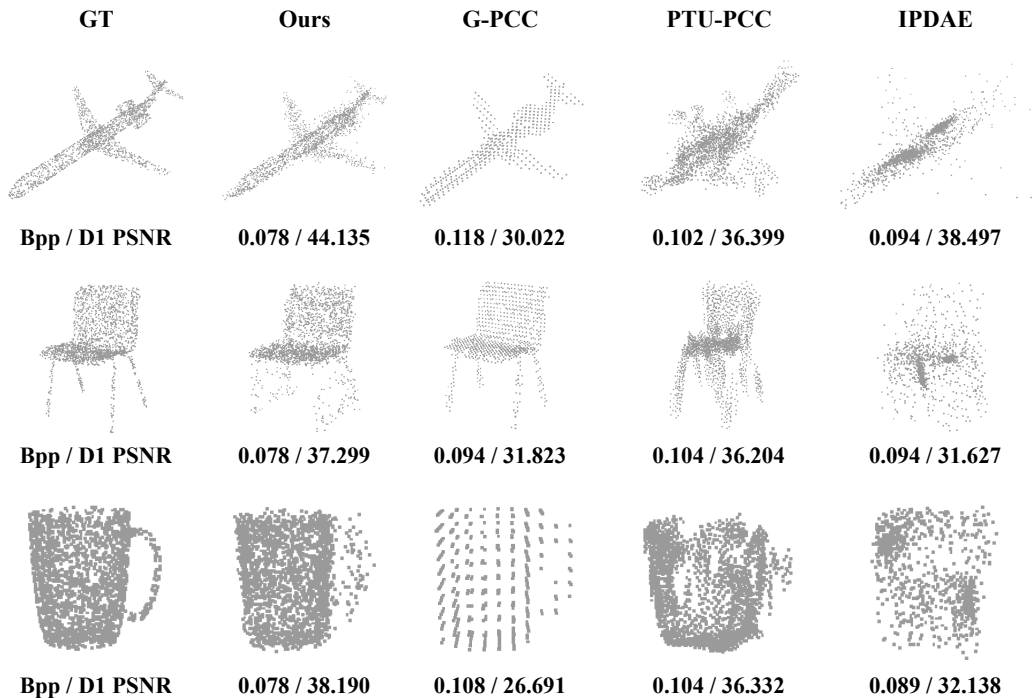


Figure 4: Subjective quality comparison. Example point clouds are selected from the ShapeNet dataset, each with 2k points.

Table 2: Ablation study of the proposed method. The original Diff-PCC serves as the anchor.

E_l backbone	E_h backbone	$\mathcal{D}_{CD}(x_0, \bar{x}_0)$	BD-PSNR (dB)	BD-Rate (%)
✓	✗	✓	-0.397	-20.637
✗	✓	✓	-2.276	-16.523
✓	✓	✗	-0.132	+4.658

232 5 Limitations

233 Although our method has achieved advanced rate distortion performance and excellent visual re-
 234 construction results, there are several limitations that warrant discussion. Firstly, the encoding and
 235 decoding time are relatively long, which could potentially be improved by the acceleration techniques
 236 employed in several explorations [18, 19]. Secondly, the model is currently limited to compressing
 237 small-scale point clouds, and further research is required to enhance its capability to handle large-scale
 238 instances.

239 6 Conclusion

240 We propose a diffusion-based point cloud compression method, dubbed Diff-PCC, to leverage the
 241 expressive power of the diffusion model for generative and aesthetically superior decoding. We
 242 introduce a dual-space latent representation to enhance the representation ability of the conventional
 243 Gaussian priors in VAEs, enabling the Diff-PCC to extract expressive shape latents and facilitate
 244 the following diffusion-based decoding process. At the decoding side, an effective diffusion-based
 245 generator produces high-quality reconstructions by considering the shape latents as guidance to
 246 stochastically denoise the noisy point clouds. The proposed method achieves state-of-the-art com-
 247 pression performance while attaining superior subjective quality. Future works may include reducing
 248 the coding complexity and extending to large-scale point cloud instances.

249 References

- 250 [1] Johannes Ballé, Valero Laparra, and Eero P Simoncelli. End-to-end optimized image compression. *arXiv preprint arXiv:1611.01704*, 2016.
251
- 252 [2] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston. Variational image
253 compression with a scale hyperprior. *arXiv preprint arXiv:1802.01436*, 2018.
- 254 [3] Francesco Paolo Casale, Adrian Dalca, Luca Saglietti, Jennifer Listgarten, and Nicolo Fusi. Gaussian
255 process prior variational autoencoders. *Advances in neural information processing systems*, 31, 2018.
- 256 [4] Ruihang Chu, Enze Xie, Shentong Mo, Zhenguo Li, Matthias Nießner, Chi-Wing Fu, and Jiaya Jia.
257 Diffcomplete: Diffusion-based generative 3d shape completion. *Advances in Neural Information Processing
258 Systems*, 36, 2024.
- 259 [5] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision:
260 A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(9):10850–10869, 2023.
- 261 [6] Bin Dai and David Wipf. Diagnosing and enhancing vae models. *arXiv preprint arXiv:1903.05789*, 2019.
- 262 [7] Kamak Ebadi, Lukas Bernreiter, Harel Biggie, Gavin Catt, Yun Chang, Arghya Chatterjee, Christopher E
263 Denniston, Simon-Pierre Deschênes, Kyle Harlow, Shehryar Khattak, et al. Present and future of slam in
264 extreme environments: The darpa sub challenge. *IEEE Transactions on Robotics*, 2023.
- 265 [8] Lili Fan, Junhao Wang, Yuanmeng Chang, Yuke Li, Yutong Wang, and Dongpu Cao. 4d mmwave radar for
266 autonomous driving perception: a comprehensive survey. *IEEE Transactions on Intelligent Vehicles*, 2024.
- 267 [9] Chunyang Fu, Ge Li, Rui Song, Wei Gao, and Shan Liu. Octattention: Octree-based large-scale contexts
268 model for point cloud compression. In *Proceedings of the AAAI conference on artificial intelligence*,
269 volume 36, pages 625–633, 2022.
- 270 [10] Xiaoran Hao and Patrick Shafto. Coupled variational autoencoder. *arXiv preprint arXiv:2306.02565*, 2023.
- 271 [11] Yun He, Xinlin Ren, Danhang Tang, Yinda Zhang, Xiangyang Xue, and Yanwei Fu. Density-preserving
272 deep point cloud compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
273 Pattern Recognition*, pages 2333–2342, 2022.
- 274 [12] Tianxin Huang, Jiangning Zhang, Jun Chen, Zhonggan Ding, Ying Tai, Zhenyu Zhang, Chengjie Wang,
275 and Yong Liu. 3qnet: 3d point cloud geometry quantization compression network. *ACM Transactions on
276 Graphics (TOG)*, 41(6):1–13, 2022.
- 277 [13] Zixuan Huang, Justin Johnson, Shoubhik Debnath, James M Rehg, and Chao-Yuan Wu. Pointinfinity:
278 Resolution-invariant point diffusion models. *arXiv preprint arXiv:2404.03566*, 2024.
- 279 [14] Yiqi Jin, Ziyu Zhu, Tongda Xu, Yuhuan Lin, and Yan Wang. Ecm-opcc: Efficient context model for
280 octree-based point cloud compression. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics,
281 Speech and Signal Processing (ICASSP)*, pages 7985–7989, 2024.
- 282 [15] Weonyoung Joo, Wonsung Lee, Sungrae Park, and Il-Chul Moon. Dirichlet variational autoencoder. *Pattern
283 Recognition*, 107:107514, 2020.
- 284 [16] M Krithika Alias AnbuDevi and K Suganthi. Review of semantic segmentation of medical images using
285 modified architectures of unet. *Diagnostics*, 12(12):3064, 2022.
- 286 [17] Jin Sub Lee, Jisun Kim, and Philip M Kim. Score-based generative modeling for de novo protein design.
287 *Nature Computational Science*, 3(5):382–392, 2023.
- 288 [18] Xiuyu Li, Yijiang Liu, Long Lian, Huanrui Yang, Zhen Dong, Daniel Kang, Shanghang Zhang, and
289 Kurt Keutzer. Q-diffusion: Quantizing diffusion models. In *Proceedings of the IEEE/CVF International
290 Conference on Computer Vision*, pages 17535–17545, 2023.
- 291 [19] Qingguo Liu, Chenyi Zhuang, Pan Gao, and Jie Qin. Cdformer: When degradation prediction embraces
292 diffusion model for blind image super-resolution. *arXiv preprint arXiv:2405.07648*, 2024.
- 293 [20] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of
294 the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021.
- 295 [21] Zhaoyang Lyu, Zhifeng Kong, Xudong Xu, Liang Pan, and Dahua Lin. A conditional point diffusion-
296 refinement paradigm for 3d point cloud completion. *ArXiv*, abs/2112.03530, 2021.

- 297 [22] David Minnen, Johannes Ballé, and George D Toderici. Joint autoregressive and hierarchical priors for
298 learned image compression. *Advances in neural information processing systems*, 31, 2018.
- 299 [23] Shentong Mo, Enze Xie, Ruihang Chu, Lanqing Hong, Matthias Niessner, and Zhenguo Li. Dit-3d:
300 Exploring plain diffusion transformers for 3d shape generation. *Advances in Neural Information Processing*
301 *Systems*, 36, 2024.
- 302 [24] Dat Thanh Nguyen and André Kaup. Lossless point cloud geometry and attribute compression using a
303 learned conditional probability model. *IEEE Transactions on Circuits and Systems for Video Technology*,
304 2023.
- 305 [25] Francesco Pezone, Osman Musa, Giuseppe Caire, and Sergio Barbarossa. Semantic-preserving image
306 coding based on conditional diffusion models. In *ICASSP 2024-2024 IEEE International Conference on*
307 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 13501–13505. IEEE, 2024.
- 308 [26] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d
309 classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern*
310 *recognition*, pages 652–660, 2017.
- 311 [27] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature
312 learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017.
- 313 [28] Guocheng Qian, Yuchen Li, Houwen Peng, Jinjie Mai, Hasan Hammoud, Mohamed Elhoseiny, and Bernard
314 Ghanem. Pointnext: Revisiting pointnet++ with improved training and scaling strategies.
- 315 [29] Ruwen Schnabel and Reinhard Klein. Octree-based point-cloud compression. *PBG@ SIGGRAPH*,
316 3:111–121, 2006.
- 317 [30] Rui Song, Chunyang Fu, Shan Liu, and Ge Li. Efficient hierarchical entropy model for learned point cloud
318 compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
319 pages 14368–14377, 2023.
- 320 [31] Rui Song, Chunyang Fu, Shan Liu, and Ge Li. Efficient hierarchical entropy model for learned point cloud
321 compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
322 pages 14368–14377, 2023.
- 323 [32] Yu Takagi and Shinji Nishimoto. High-resolution image reconstruction with latent diffusion models from
324 human brain activity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
325 *Recognition*, pages 14453–14463, 2023.
- 326 [33] Lucas Theis, Tim Salimans, Matthew D Hoffman, and Fabian Mentzer. Lossy compression with gaussian
327 diffusion. *arXiv preprint arXiv:2206.08889*, 2022.
- 328 [34] Anwaar Ulhaq, Naveed Akhtar, and Ganna Pogrebna. Efficient diffusion models for vision: A survey.
329 *arXiv preprint arXiv:2210.09292*, 2022.
- 330 [35] Juho-Pekka Virtanen, Sylvie Daniel, Tuomas Turppa, Lingli Zhu, Arttu Julin, Hannu Hyypä, and Juha
331 Hyypä. Interactive dense point clouds in a game engine. *ISPRS Journal of Photogrammetry and Remote*
332 *Sensing*, 163:375–389, 2020.
- 333 [36] Jianqiang Wang, Dandan Ding, Zhu Li, and Zhan Ma. Multiscale point cloud geometry compression. In
334 *2021 Data Compression Conference (DCC)*, pages 73–82. IEEE, 2021.
- 335 [37] Jianqiang Wang, Dandan Ding, and Zhan Ma. Lossless point cloud attribute compression using cross-scale,
336 cross-group, and cross-color prediction. In *2023 Data Compression Conference (DCC)*, pages 228–237.
337 IEEE, 2023.
- 338 [38] Jianqiang Wang and Zhan Ma. Sparse tensor-based point cloud attribute compression. In *2022 IEEE*
339 *5th International Conference on Multimedia Information Processing and Retrieval (MIPR)*, pages 59–64.
340 IEEE, 2022.
- 341 [39] MPEG 3D Graphics WG 7 and Haptics Coding. G-pcc 2nd edition codec description. *ISO/IEC JTC 1/SC*
342 *29/WG 7*, 2023.
- 343 [40] MPEG 3D Graphics Coding WG 7. V-pcc codec description. *ISO/IEC JTC 1/SC 29/WG 7*, 2020.
- 344 [41] Yankun Wu, Yuta Nakashima, and Noa Garcia. Not only generative art: Stable diffusion for content-style
345 disentanglement in art analysis. In *Proceedings of the 2023 ACM International conference on multimedia*
346 *retrieval*, pages 199–208, 2023.

- 347 [42] Ruixiang Xue, Jiaxin Li, Tong Chen, Dandan Ding, Xun Cao, and Zhan Ma. Neri: Implicit neural
348 representation of lidar point cloud using range image sequence. In *ICASSP 2024-2024 IEEE International*
349 *Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8020–8024. IEEE, 2024.
- 350 [43] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow:
351 3d point cloud generation with continuous normalizing flows. *arXiv*, 2019.
- 352 [44] Ruihan Yang and Stephan Mandt. Lossy image compression with conditional diffusion models. *Advances*
353 *in Neural Information Processing Systems*, 36, 2024.
- 354 [45] Kang You, Pan Gao, and Qing Li. Ipdac: Improved patch-based deep autoencoder for lossy point cloud
355 geometry compression. In *Proceedings of the 1st International Workshop on Advances in Point Cloud*
356 *Compression, Processing and Analysis*, pages 1–10, 2022.
- 357 [46] Kang You, Kai Liu, Li Yu, Pan Gao, and Dandan Ding. Pointsoup: High-performance and ex-
358 tremely low-decoding-latency learned geometry codec for large-scale point cloud scenes. *arXiv preprint*
359 *arXiv:2404.13550*, 2024.
- 360 [47] Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, and Karsten Kreis.
361 Lion: Latent point diffusion models for 3d shape generation. In *Advances in Neural Information Processing*
362 *Systems (NeurIPS)*, 2022.
- 363 [48] Chenshuang Zhang, Chaoning Zhang, Mengchun Zhang, and In So Kweon. Text-to-image diffusion model
364 in generative ai: A survey. *arXiv preprint arXiv:2303.07909*, 2023.
- 365 [49] Junteng Zhang, Tong Chen, Dandan Ding, and Zhan Ma. Yoga: Yet another geometry-based point cloud
366 compressor. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 9070–9081,
367 2023.
- 368 [50] Junteng Zhang, Gexin Liu, Dandan Ding, and Zhan Ma. Transformer and upsampling-based point cloud
369 compression. In *Proceedings of the 1st International Workshop on Advances in Point Cloud Compression,*
370 *Processing and Analysis*, pages 33–39, 2022.
- 371 [51] Junteng Zhang, Jianqiang Wang, Dandan Ding, and Zhan Ma. Scalable point cloud attribute compression.
372 *IEEE Transactions on Multimedia*, 2023.
- 373 [52] Junzhe Zhang, Tong Chen, Dandan Ding, and Zhan Ma. G-pcc++: Enhanced geometry-based point cloud
374 compression. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 1352–1363,
375 2023.
- 376 [53] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion
377 models.
- 378 [54] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion
379 models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847,
380 2023.
- 381 [55] Renrui Zhang, Liuhui Wang, Ziyu Guo, Yali Wang, Peng Gao, Hongsheng Li, and Jianbo Shi. Parameter
382 is not all you need: Starting from non-parametric networks for 3d point cloud analysis. *arXiv preprint*
383 *arXiv:2303.08134*, 2023.
- 384 [56] Shengjia Zhao, Jiaming Song, and Stefano Ermon. Towards deeper understanding of variational autoen-
385 coding models. *arXiv preprint arXiv:1702.08658*, 2017.
- 386 [57] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion.
387 In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5826–5835,
388 October 2021.
- 389 [58] Yuanzhi Zhu, Kai Zhang, Jingyun Liang, Jiezhong Cao, Bihan Wen, Radu Timofte, and Luc Van Gool.
390 Denoising diffusion models for plug-and-play image restoration. In *Proceedings of the IEEE/CVF*
391 *Conference on Computer Vision and Pattern Recognition*, pages 1219–1229, 2023.

392 NeurIPS Paper Checklist

393 1. Claims

394 Question: Do the main claims made in the abstract and introduction accurately reflect the
395 paper's contributions and scope?

396 Answer: [Yes]

397 Justification: Claims are clearly stated in abstract and introduction (Sec. 1). The experimental
398 results (Sec. 4) match with these claims and reflect how much the results can be expected to
399 generalize to other settings.

400 Guidelines:

- 401 • The answer NA means that the abstract and introduction do not include the claims
402 made in the paper.
- 403 • The abstract and/or introduction should clearly state the claims made, including the
404 contributions made in the paper and important assumptions and limitations. A No or
405 NA answer to this question will not be perceived well by the reviewers.
- 406 • The claims made should match theoretical and experimental results, and reflect how
407 much the results can be expected to generalize to other settings.
- 408 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
409 are not attained by the paper.

410 2. Limitations

411 Question: Does the paper discuss the limitations of the work performed by the authors?

412 Answer: [Yes]

413 Justification: Limitations are discussed in Sec. 5.

414 Guidelines:

- 415 • The answer NA means that the paper has no limitation while the answer No means that
416 the paper has limitations, but those are not discussed in the paper.
- 417 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 418 • The paper should point out any strong assumptions and how robust the results are to
419 violations of these assumptions (e.g., independence assumptions, noiseless settings,
420 model well-specification, asymptotic approximations only holding locally). The authors
421 should reflect on how these assumptions might be violated in practice and what the
422 implications would be.
- 423 • The authors should reflect on the scope of the claims made, e.g., if the approach was
424 only tested on a few datasets or with a few runs. In general, empirical results often
425 depend on implicit assumptions, which should be articulated.
- 426 • The authors should reflect on the factors that influence the performance of the approach.
427 For example, a facial recognition algorithm may perform poorly when image resolution
428 is low or images are taken in low lighting. Or a speech-to-text system might not be
429 used reliably to provide closed captions for online lectures because it fails to handle
430 technical jargon.
- 431 • The authors should discuss the computational efficiency of the proposed algorithms
432 and how they scale with dataset size.
- 433 • If applicable, the authors should discuss possible limitations of their approach to
434 address problems of privacy and fairness.
- 435 • While the authors might fear that complete honesty about limitations might be used by
436 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
437 limitations that aren't acknowledged in the paper. The authors should use their best
438 judgment and recognize that individual actions in favor of transparency play an impor-
439 tant role in developing norms that preserve the integrity of the community. Reviewers
440 will be specifically instructed to not penalize honesty concerning limitations.

441 3. Theory Assumptions and Proofs

442 Question: For each theoretical result, does the paper provide the full set of assumptions and
443 a complete (and correct) proof?

444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496

Answer: [NA]

Justification: This paper does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The proposed architecture is fully described in Sec. 3, detailed instructions for replication is provided in the experimental setup section (Sec. 4.1).

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

497 Question: Does the paper provide open access to the data and code, with sufficient instruc-
498 tions to faithfully reproduce the main experimental results, as described in supplemental
499 material?

500 Answer: [No]

501 Justification: Source code will be made publicly available once paper is accepted.

502 Guidelines:

- 503 • The answer NA means that paper does not include experiments requiring code.
- 504 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
505 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 506 • While we encourage the release of code and data, we understand that this might not be
507 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
508 including code, unless this is central to the contribution (e.g., for a new open-source
509 benchmark).
- 510 • The instructions should contain the exact command and environment needed to run to
511 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
512 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 513 • The authors should provide instructions on data access and preparation, including how
514 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 515 • The authors should provide scripts to reproduce all experimental results for the new
516 proposed method and baselines. If only a subset of experiments are reproducible, they
517 should state which ones are omitted from the script and why.
- 518 • At submission time, to preserve anonymity, the authors should release anonymized
519 versions (if applicable).
- 520 • Providing as much information as possible in supplemental material (appended to the
521 paper) is recommended, but including URLs to data and code is permitted.

522 6. Experimental Setting/Details

523 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
524 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
525 results?

526 Answer: [Yes]

527 Justification: Experimental setting and details are fully disclosed in Sec. 4.

528 Guidelines:

- 529 • The answer NA means that the paper does not include experiments.
- 530 • The experimental setting should be presented in the core of the paper to a level of detail
531 that is necessary to appreciate the results and make sense of them.
- 532 • The full details can be provided either with the code, in appendix, or as supplemental
533 material.

534 7. Experiment Statistical Significance

535 Question: Does the paper report error bars suitably and correctly defined or other appropriate
536 information about the statistical significance of the experiments?

537 Answer: [No]

538 Justification: Error bars are not reported due to the specificity of the compression task. The
539 rate-distortion curve (Fig. 3) and Bjontegaard metric (Tab. 1) could be convincing enough.

540 Guidelines:

- 541 • The answer NA means that the paper does not include experiments.
- 542 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
543 dence intervals, or statistical significance tests, at least for the experiments that support
544 the main claims of the paper.
- 545 • The factors of variability that the error bars are capturing should be clearly stated (for
546 example, train/test split, initialization, random drawing of some parameter, or overall
547 run with given experimental conditions).

- 548 • The method for calculating the error bars should be explained (closed form formula,
549 call to a library function, bootstrap, etc.)
- 550 • The assumptions made should be given (e.g., Normally distributed errors).
- 551 • It should be clear whether the error bar is the standard deviation or the standard error
552 of the mean.
- 553 • It is OK to report 1-sigma error bars, but one should state it. The authors should
554 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
555 of Normality of errors is not verified.
- 556 • For asymmetric distributions, the authors should be careful not to show in tables or
557 figures symmetric error bars that would yield results that are out of range (e.g. negative
558 error rates).
- 559 • If error bars are reported in tables or plots, The authors should explain in the text how
560 they were calculated and reference the corresponding figures or tables in the text.

561 8. Experiments Compute Resources

562 Question: For each experiment, does the paper provide sufficient information on the com-
563 puter resources (type of compute workers, memory, time of execution) needed to reproduce
564 the experiments?

565 Answer: [Yes]

566 Justification: Sufficient information on the computer resources is disclosed in the experiment
567 setting (Sec. 4.1).

568 Guidelines:

- 569 • The answer NA means that the paper does not include experiments.
- 570 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
571 or cloud provider, including relevant memory and storage.
- 572 • The paper should provide the amount of compute required for each of the individual
573 experimental runs as well as estimate the total compute.
- 574 • The paper should disclose whether the full research project required more compute
575 than the experiments reported in the paper (e.g., preliminary or failed experiments that
576 didn't make it into the paper).

577 9. Code Of Ethics

578 Question: Does the research conducted in the paper conform, in every respect, with the
579 NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

580 Answer: [Yes]

581 Justification: This research does not involve human subjects or participants. This paper
582 conform with the Code of Ethics in every respect.

583 Guidelines:

- 584 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 585 • If the authors answer No, they should explain the special circumstances that require a
586 deviation from the Code of Ethics.
- 587 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
588 eration due to laws or regulations in their jurisdiction).

589 10. Broader Impacts

590 Question: Does the paper discuss both potential positive societal impacts and negative
591 societal impacts of the work performed?

592 Answer: [NA]

593 Justification: There is no societal impact of the work. The proposed method is limited to
594 compression and reconstruction and cannot be used to generate deepfakes or disinformation.

595 Guidelines:

- 596 • The answer NA means that there is no societal impact of the work performed.
- 597 • If the authors answer NA or No, they should explain why their work has no societal
598 impact or why the paper does not address societal impact.

- 599
- 600
- 601
- 602
- 603
- 604
- 605
- 606
- 607
- 608
- 609
- 610
- 611
- 612
- 613
- 614
- 615
- 616
- 617
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

618 11. Safeguards

619 Question: Does the paper describe safeguards that have been put in place for responsible
620 release of data or models that have a high risk for misuse (e.g., pretrained language models,
621 image generators, or scraped datasets)?

622 Answer: [NA]

623 Justification: This paper poses no safeguard risks.

624 Guidelines:

- 625
- 626
- 627
- 628
- 629
- 630
- 631
- 632
- 633
- 634
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

635 12. Licenses for existing assets

636 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
637 the paper, properly credited and are the license and terms of use explicitly mentioned and
638 properly respected?

639 Answer: [Yes]

640 Justification: This paper follows the license of the datasets used. Original papers are properly
641 cited.

642 Guidelines:

- 643
- 644
- 645
- 646
- 647
- 648
- 649
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- 650
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

658 13. **New Assets**

659 Question: Are new assets introduced in the paper well documented and is the documentation
660 provided alongside the assets?

661 Answer: [NA]

662 Justification: This paper does not release new assets.

663 Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

672 14. **Crowdsourcing and Research with Human Subjects**

673 Question: For crowdsourcing experiments and research with human subjects, does the paper
674 include the full text of instructions given to participants and screenshots, if applicable, as
675 well as details about compensation (if any)?

676 Answer: [NA]

677 Justification: This paper does not involve crowdsourcing nor research with human subjects.

678 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

687 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

689 Question: Does the paper describe potential risks incurred by study participants, whether
690 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
691 approvals (or an equivalent approval/review based on the requirements of your country or
692 institution) were obtained?

693 Answer: [NA]

694 Justification: This paper does not involve crowdsourcing nor research with human subjects.

695 Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

701
702
703
704
705

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.