

ContextGS: Compact 3D Gaussian Splatting with Anchor Level Context Model

Yufei Wang¹ Zhihao Li¹ Lanqing Guo¹ Wenhan Yang² Alex C. Kot¹ Bihan Wen¹

¹ Nanyang Technological University, Singapore

² PengCheng Laboratory, China

{yufei001, zhihao.li, lanqing.guo, eackot, bihan.wen}@ntu.edu.sg yangwh@pcl.ac.cn

Homepage: <https://github.com/wyf0912/ContextGS>

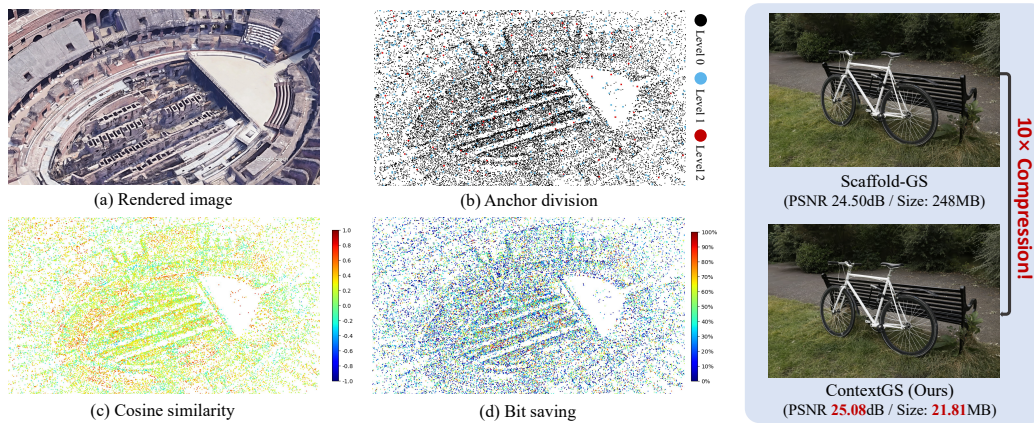


Figure 1: An illustration of the necessity of using autoregressive model in the anchor level. While Scaffold-GS [20] greatly reduces the spatial redundancy among adjacent 3D Gaussians by grouping them and introducing a new data structure *anchor* to capture their common features, spatial redundancy still exists among anchors. Our method, **ContextGS**, first proposes to reduce the spatial redundancy among anchors using an autoregressive model. We divide anchors into levels as shown in Fig. (b) and the anchors from coarser levels are used to predict anchors in finer levels, *i.e.*, \bullet predicts \bullet then \bullet together predict \bullet . Fig. (c) verifies the spatial redundancy by calculating the cosine similarity between anchors in level 0 and their context anchors in levels 1 and 2. Fig. (d) displays the bit savings using the proposed anchor-level context model evaluated on our entropy coding based strong baseline built on Scaffold-GS [20]. Compared with Scaffold-GS, we achieve better rendering qualities, faster rendering speed, and great size reduction of up to 15 times averaged over all datasets we used.

Abstract

Recently, 3D Gaussian Splatting (3DGS) has become a promising framework for novel view synthesis, offering fast rendering speeds and high fidelity. However, the large number of Gaussians and their associated attributes require effective compression techniques. Existing methods primarily compress 3D Gaussians individually and independently, *i.e.*, coding all the 3D Gaussians at the same time, with little design for their interactions and spatial dependence. Inspired by the effectiveness of the context model in image compression, we propose the first autoregressive model at the anchor level for 3DGS compression in this work. We divide anchors into different levels and the anchors that are not coded yet can be predicted based on the already coded ones in all the coarser levels, leading to more

accurate modeling and higher coding efficiency. To further improve the efficiency of entropy coding, *e.g.*, to code the coarsest level with no already coded anchors, we propose to introduce a low-dimensional quantized feature as the hyperprior for each anchor, which can be effectively compressed. Our work pioneers the context model in the anchor level for 3DGS representation, yielding an impressive size reduction of over 100 times compared to vanilla 3DGS and 15 times compared to the most recent state-of-the-art work Scaffold-GS, while achieving comparable or even higher rendering quality.

1 Introduction

Over the past few years, novel view synthetic has rapidly progressed. As a representative work, Neural Radiance Field (NeRF) [22] uses a Multilayer Perceptron (MLP) to predict the attributes of queried points in the 3D scene. While good rendering qualities are achieved, the dense querying process results in slow rendering, which greatly hinders their applications in practical scenarios. Significant efforts have been made to enhance training and rendering speeds, achieving notable progress through various techniques, such as factorization [4, 8, 10, 13] and hash grids [24, 12]. However, they still face challenges in the real-time rendering of large-scale scenes due to the intrinsic limitations of volumetric sampling. Recently, 3D Gaussian Splatting (3DGS) [15] has achieved state-of-the-art (SOTA) rendering quality and speed. As an emerging alternative strategy for representing 3D scenes, 3DGS represents a 3D scene using a set of 3D Gaussians initiated from Structure-from-Motion (SfM) with learnable attributes such as color, shape, and opacity. The 2D images can be effectively rendered using differentiable rasterization and end-to-end training is enabled. Meanwhile, benefiting from efficient CUDA implementation, real-time rendering is achieved.

Despite its success, 3DGS still encounters limitations in storage efficiency. Representing large scenes requires millions of 3D Gaussian points, which demand several GBs of storage, *e.g.*, an average of 1.6 GB for each scene in the BungeeNerf [32] dataset. The huge storage overhead greatly hinders the applications of 3DGS [15], thus an efficient compression technique is required. However, the unorganized and sparse nature of these 3D Gaussians makes it highly challenging to effectively reduce data redundancy. To address this issue, various techniques have been proposed to enhance the storage efficiency of 3D Gaussian models. For example, [7, 18, 25, 26] proposed to discretize continuous attributes of 3D Gaussians to a cluster of attributes stored in the codebooks; [7, 18] proposed to prune neural Gaussians with little effect. Entropy coding is also used to reduce the storage overhead by further encoding neural Gaussian features into bitstream [11, 5, 18, 23]. Although space utilization has greatly improved, they focus more on individually compressing each Gaussian point and do not well explore the relationship and reduce the spatial redundancy among neural Gaussians. To further reduce the spatial redundancy, most recently, [20] proposed to divide anchors into voxels and introduced an anchor feature for each voxel to grasp the common attributes of neural Gaussians in the voxel, *i.e.*, the neural Gaussians are predicted by the anchor features. While the spatial dependency has been significantly reduced, as shown in Fig. 1 (c), the similarity among anchors remains high in certain areas, indicating that spatial redundancy still exists.

To further enhance the coding efficiency of 3D scenes, we propose a novel framework named *ContextGS* for 3DGS compression. Inspired by the effectiveness of context models [31] in image compression [21], we introduce an autoregressive model at the anchor level into 3DGS. Specifically, building on Scaffold-GS [20], we divide anchors into hierarchical levels and encode them progressively. Coarser level anchors are encoded first, and their decoded values are used to predict the distribution of nearby anchors at finer levels. This approach leverages spatial dependencies among adjacent anchors, allowing already decoded anchors to better predict the distribution of subsequent anchors, leading to significant improvements in coding efficiency. Additionally, anchors decoded at coarser levels can be directly used in the final fine-grained level, reducing storage overhead. To further enhance coding efficiency, especially for encoding the coarsest level anchors without already decoded ones, we employ a quantized hyperprior feature as an additional prior for each anchor. It is worth noting that the proposed method can also support vanilla 3DGS renderers by simply removing the view-dependent features from the anchor feature. Our contributions can be summarized as follows:

- We propose the first context model for 3DGS at the anchor level. By predicting the properties of anchors that are not coded yet given already coded ones, we greatly eliminate the spatial redundancy among anchors.

- We propose a unified compressing framework with the factorized prior, enabling end-to-end entropy coding of anchor features. Besides, a strategy for anchor layering is proposed, which allows already decoded anchors to quickly locate adjacent anchors that are to be decoded. Meanwhile, the proposed method avoids redundant coding of anchors by the proposed anchor forward.
- The experimental results on real-world datasets demonstrate the effectiveness of the proposed method compared with SOTA and concurrent works. On average across all datasets, our model achieves a compression ratio of $15\times$ compared to the Scaffold-GS model we used as the backbone and $100\times$ compared to the standard 3DGS model, while maintaining comparable or even enhanced fidelity.

2 Related works

2.1 Neural radiance field and 3D Gaussian splatting

Early 3D scene modeling often employs the Neural Radiance Field (NeRF) [22] as a global approximator for 3D scene appearance and geometry. These approaches [2, 29, 30] use a multi-layer perceptron (MLP) to implicitly represent the 3D scene by predicting attributes of queried points. However, the dense querying process results in extremely slow rendering. Various methods have been developed to speed up the rendering process significantly [6, 9, 27], such as plane factorization-based techniques like K-Planes [8] and TensorRF [4], and the use of hash grid features in InstantNGP [24]. While these methods enable high-quality rendering with a much smaller MLP compared to the vanilla NeRF, rendering a single pixel still requires numerous queries. This can lead to increased storage requirements for the grid-based features and difficulties in efficiently rendering empty space or large-scale scenes. To achieve real-time and efficient rendering while maintaining high fidelity, 3DGS [15] introduces an innovative approach by representing the scene explicitly with numerous learnable 3D Gaussians. By employing differentiable splatting and tile-based rasterization [17], 3DGS [15] optimizes these Gaussians during training in an end-to-end manner.

2.2 Deep compression

Despite the effectiveness of 3DGS [15] in rendering speed and fidelity, the large number of Gaussians and their associated attributes result in significant storage overhead. Many techniques are proposed to reduce the storage requirements of 3DGS. For example, [7, 18] proposes to prune neural Gaussians (aka 3D Gaussians) with minimal impact. [7, 18, 25, 26] propose to utilize codebooks to cluster Gaussian parameters. Entropy coding is also used in [11, 5, 18, 23] to encode neural Gaussians into bit streams by modeling their distributions. While remarkable performances are achieved, they mainly focus on improving the efficiency of a single neural Gaussian and neglect the spatial redundancy among neighbor neural Gaussians. Most recently, Scaffold-GS [20] proposes to introduce an anchor level to capture common features of nearby neural Gaussians in the same voxel, and successive work [5] demonstrates its effectiveness by further introducing hash-feature as a prior for entropy coding. However, [5] codes all the anchors at the same time, and its spatial redundancy can be further reduced. In the image compression task, an important category of methods to improve the coding efficiency is the context model [21, 31], which greatly reduces the spatial redundancy by predicting the distribution of latent pixels based on an already coded one. Inspired by the context models used in compression, we propose to encode the anchor features in an autoregressive way, *i.e.*, predict the anchor points from already coded ones at coarser levels. As far as we know, we are the first to reduce the storage redundancy of 3DGS using a context model at the anchor level.

3 Preliminary

3DGS [15] utilizes a collection of anisotropic 3D neural Gaussians to depict the scene so that the scene can be efficiently rendered using a tile-based rasterization technique. Beginning from a set of Structure-from-Motion (SfM) points, each Gaussian point is represented as follows

$$G(\mathbf{x}) = e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}, \quad (1)$$

where \mathbf{x} is the coordinates in the 3D scene, $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are the mean position and covariance matrix of the Gaussian point, respectively. To ensure the positive semi-definite of $\boldsymbol{\Sigma}$, $\boldsymbol{\Sigma}$ is represented as

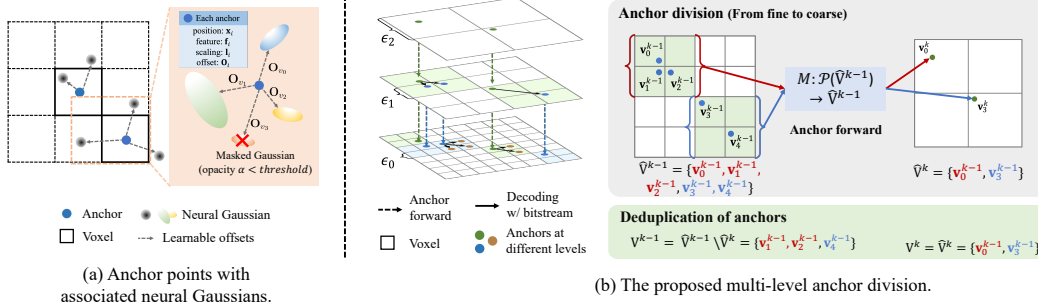


Figure 2: **(a)**: An illustration of the data structure we used following Scaffold-GS [19], where anchor points are used to extract common features of their associated neural Gaussians. **(b)**: The proposed multi-level division of anchor points. The decoded anchors from higher (coarser) levels are directly forwarded to the lower (finer) level to avoid duplicate storage. Besides, taking decompression as an example, the already decoded anchors are used to predict anchors that are not decompressed yet, which greatly reduces the spatial redundancy among adjacent anchors. (Best zoom in for details.)

$\Sigma = \mathbf{R}\mathbf{S}\mathbf{S}^T\mathbf{R}^T$, where \mathbf{R} and \mathbf{S} are scaling and rotation matrixes, respectively. Besides, each neural Gaussian has the attributes of opacity $\alpha \in \mathbb{R}^1$ and view-dependent color $\mathbf{c} \in \mathbb{R}^3$ modeled by spherical harmonic [34]. All the attributes, *i.e.*, $[\boldsymbol{\mu}, \mathbf{R}, \mathbf{S}, \alpha, \mathbf{c}]$, in neural Gaussian points are learnable and optimized by the reconstruction loss of images rendered by the tile-based rasterization.

Scaffold-GS [20]. While representing scenes with neural Gaussians greatly accelerates the rendering speed, the large amount of 3D Gaussians leads to significant storage overhead. To reduce the redundancy among adjunct 3D Gaussians, the most recent work, Scaffold-GS [19], proposes to introduce *anchor points* to capture common attributes of local 3D Gaussians as shown in Fig. 2 (a). Specifically, the *anchor points* are initialized from neural Gaussians by voxelizing the 3D scenes. Each anchor point has a context feature $\mathbf{f} \in \mathbb{R}^{32}$, a location $\mathbf{x} \in \mathbb{R}^3$, a scaling factor $\mathbf{l} \in \mathbb{R}^3$ and k learnable offset $\mathbf{O} \in \mathbb{R}^{k \times 3}$. Given a camera at \mathbf{x}_c , anchor points are used to predict the view-dependent neural Gaussians in their corresponding voxels as follows,

$$\{\mathbf{c}^i, \mathbf{r}^i, \mathbf{s}^i, \alpha^i\}_{i=0}^k = F(\mathbf{f}, \boldsymbol{\sigma}_c, \vec{\mathbf{d}}_c) \quad (2)$$

where $\boldsymbol{\sigma}_c = \|\mathbf{x} - \mathbf{x}_c\|_2$, $\vec{\mathbf{d}}_c = \frac{\mathbf{x} - \mathbf{x}_c}{\|\mathbf{x} - \mathbf{x}_c\|_2}$, the superscript i represents the index of neural Gaussian in the voxel, $\mathbf{s}^i, \mathbf{c}^i \in \mathbb{R}^3$ are the scaling and color respectively, and $\mathbf{r}^i \in \mathbb{R}^4$ is the quaternion for rotation. The positions of neural Gaussians are then calculated as

$$\{\boldsymbol{\mu}^0, \dots, \boldsymbol{\mu}^{k-1}\} = \mathbf{x} + \{\mathbf{O}^0, \dots, \mathbf{O}^{k-1}\} \cdot \mathbf{l}, \quad (3)$$

where \mathbf{x} is the learnable positions of the anchor and \mathbf{l} is the base scaling of its associated neural Gaussians. After decoding the properties of neural Gaussians from anchor points, other processes are the same as the 3DGS [15]. By predicting the properties of neural Gaussians from the anchor features and saving the properties of anchor points only, Scaffold-GS [19] greatly eliminates the redundancy among 3D neural Gaussians and decreases the storage demand.

4 Methodology

The overall framework is shown in Fig. 3. We first introduce how to divide anchors into levels with traceable mapping relationships among adjacent levels in Sec 4.1. Based on that, we present the entropy coding in an autoregressive way in Sec 4.2, and the overall training objective in Sec 4.3.

4.1 Anchor partitioning strategy

We attempt to divide N anchors $\mathbf{V} = \{\mathbf{v}_i\}_{i=0}^N = \{(\mathbf{x}_i, \mathbf{f}_i, \mathbf{l}_i, \mathbf{O}_i)\}_{i=0}^N$ into K disjoint levels, *i.e.*, $\mathbf{V} = \mathbf{V}^0 \cup \mathbf{V}^1 \cup \dots \cup \mathbf{V}^{K-1}$ and $\mathbf{V}^i \cap \mathbf{V}^j = \emptyset, \forall i \neq j$. For each anchor set \mathbf{V}^i , it is expected to spawn over the whole scene, and be a relatively uniform downsampling of a finer set \mathbf{V}^{i-1} . Assume that we encode/decode the scene using the order from level $K-1$ to level 0 where level $K-1$ is

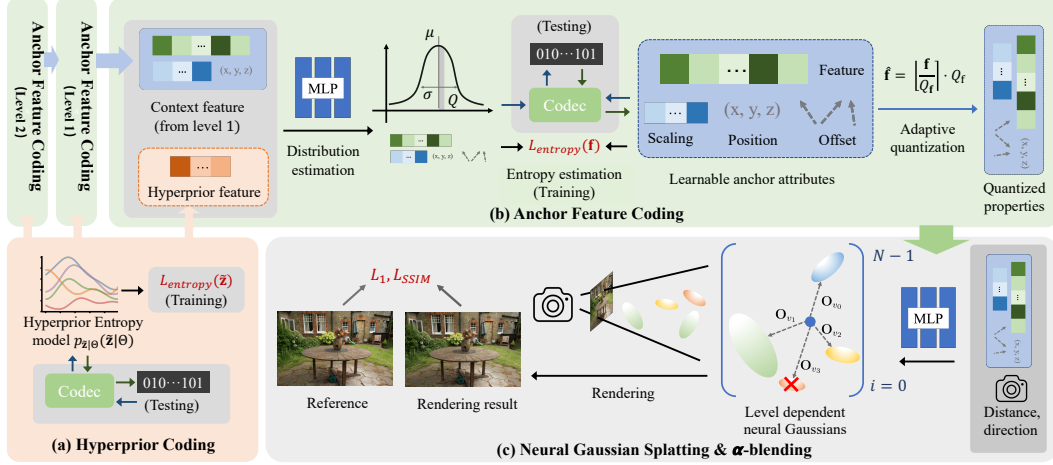


Figure 3: The overall framework of the proposed method includes three levels, *i.e.*, $K = 3$, to encode the anchors. The decoded anchors from a coarser level $i + 1$ are used to encode the anchors in level i . Besides, hyperprior features are used to predict the properties of anchors at all levels. For training, after finishing the coding of all levels, the anchor features after adaptive quantization are used to predict properties of neural Gaussians. The rendering loss is calculated and optimized together with the entropy coding loss $\mathcal{L}_{entropy}$. For testing, after we decode anchor features from the bit stream, the rendering is exactly the same with Scaffold-GS [20] without introducing overhead.

the coarsest level, we expect the mapping from \mathbf{V}^i to \mathbf{V}^{i-1} is traceable and easy to obtain. In other words, given an anchor $\mathbf{v}_i^k = (\mathbf{x}_i^k, \mathbf{f}_i^k, l_i^k, \mathbf{O}_i^k)$ from a coarser level k , we expect to quickly locate \mathbf{v}_i^k 's adjacent anchor set $\{\mathbf{v}_i^{k-1}\}_{i=0}^{N_i^{k-1}}$ in level $k - 1$, where N_i^{k-1} is the number of adjacent anchors, *i.e.*, the anchors in the same voxel of level $k - 1$ with \mathbf{v}_i^k . As such, after decoding anchors at a coarser level, we can scatter the already decoded features to to-be-processed anchors as the prior for better coding efficiency.

To achieve the requirements above, we propose to utilize a simple yet effective way to divide anchors into levels using a “bottom-up” method. As shown in Fig. 2 (b), given a set of anchors of the scene, we partition them into different fine-grained sets based on different voxel sizes ϵ_i as follows

$$\hat{\mathbf{V}}^k = \left\{ M(\{\mathbf{v}_i^{k-1} : \hat{\mathbf{x}}_i^k = \hat{\mathbf{x}}\}) : \hat{\mathbf{x}} \in \{\hat{\mathbf{x}}_i^k : i = 1, 2, \dots, |\hat{\mathbf{V}}^{k-1}|\} \right\}, \quad (4)$$

where $|\cdot|$ is the counting operation, $\hat{\mathbf{x}}_i^k$ is the anchor position after quantization using the voxel size of level k , and $M : \mathcal{P}(\hat{\mathbf{V}}^{k-1}) \rightarrow \hat{\mathbf{V}}^{k-1}$ (where \mathcal{P} is the power set) is a mapping function that selects a representative anchor to level $k - 1$ from a set of anchors $\{\mathbf{v}_i^{k-1} : \hat{\mathbf{x}}_i^k = \hat{\mathbf{x}}\}$ that has the same position after quantization. The definition of $\hat{\mathbf{x}}_i^k$ and M are as follows

$$\hat{\mathbf{x}}_i^k = \left\lfloor \frac{\mathbf{x}_i^{k-1}}{\epsilon^k} \right\rfloor \epsilon^k, \quad \mathbf{x}_i^{k-1} \in \hat{\mathbf{V}}^{k-1} \quad (5)$$

$$M(\{\mathbf{v}_i^{k-1} : \hat{\mathbf{x}}_i^k = \hat{\mathbf{x}}\}) = \mathbf{v}_j^{k-1} \text{ s.t. } j = \min\{i : \hat{\mathbf{x}}_i^k = \hat{\mathbf{x}}\},$$

where $\hat{\mathbf{V}}^0$ is initialized as the whole anchor set \mathbf{V} . We select the anchor with the minimum index $\min\{i : \hat{\mathbf{x}}_i^k = \hat{\mathbf{x}}\}$ among a set of anchors in level $k - 1$ that are in the same voxel. Besides, we filter out the repeated anchors among different levels as follows

$$\mathbf{V}^2 = \hat{\mathbf{V}}^2, \mathbf{V}^1 = \hat{\mathbf{V}}^1 \setminus \hat{\mathbf{V}}^2, \mathbf{V}^0 = \hat{\mathbf{V}}^0 \setminus \hat{\mathbf{V}}^1 \quad (6)$$

where \setminus is the set difference operation. We keep the voxel size ϵ_0 of the finest level (level 0) the same as the initial value ϵ , and set the voxel size of level i to $\epsilon_i = \kappa_i \cdot \epsilon$. Since different scenes have different initial voxel sizes and anchor point distributions, using a fixed set of voxel scaling parameters $\{\kappa_i\}_{i=1}^K$ leads to a suboptimal performance. To avoid finetuning hyper-parameters for each scene, we propose to conduct a one-time parameter search after initializing anchors. Instead of

directly setting the scale κ_i , we set a target ratio between level i and $i + 1$ and expect $\frac{|\mathbf{V}^{i+1}|}{|\mathbf{V}^i|} \approx \tau$. Since $|\mathbf{V}^i|$ decreases monotonically with κ_i , we can easily and efficiently determine the values of $\{\kappa_i\}_{i=1}^K$ using a binary search. We empirically find that the performance of models among different senses is relatively robust to the selection of τ (refer to Fig. 6).

4.2 Coding with entropy models

After dividing the anchors into multi-levels, in this section, we discuss how we use the already decoded anchors to predict ones that are not decompressed yet and how to encode attributes of anchors to improve the coding efficiency.

Context modeling in anchor levels. To encode an anchor point $\mathbf{v}_i = (\mathbf{x}_i, \mathbf{f}_i, \mathbf{l}_i, \mathbf{O}_i)$ into bitstreams efficiently using entropy coding, we need to estimate its distributions accurately. The core idea of the proposed method is to predict the properties of anchors additionally conditioned on already decompressed anchors. Taking the modeling of anchor feature \mathbf{f}_i^k from the anchor \mathbf{v}_i^k as an example, the details are as follows

$$p_{\mathbf{f}_i^k}(\mathbf{f}_i^k | \psi_i^k) = (\mathcal{N}(\boldsymbol{\mu}_i^k, \boldsymbol{\sigma}_i^k) * \mathcal{U}(-\frac{1}{2}\Delta_i^k, \frac{1}{2}\Delta_i^k))(\mathbf{f}_i^k), \quad \boldsymbol{\mu}_i^k, \boldsymbol{\sigma}_i^k, \Delta_i^k = F_{\mathbf{f}_i^k}^k(\psi_i^k), \quad (7)$$

where $F_{\mathbf{f}_i^k}^k$ is a MLP belonging to the level k , and ψ_i^k is the prior of the anchor \mathbf{v}_i^k as follows

$$\psi_i^k = \begin{cases} [\mathbf{x}_i^k] & k = K - 1 \\ [\mathbf{f}_j^{k+1}; \mathbf{l}_j^{k+1}; \mathbf{x}_i^k], & k < K - 1 \end{cases} \quad (8)$$

where $[\cdot]$ is the concatenation operation among the channel dimension, $\mathbf{f}_j^{k+1}, \mathbf{l}_j^{k+1}$ are the feature and scaling from the adjacent anchor \mathbf{v}_j^{k+1} in level $k + 1$ that are already decoded as shown in Fig. 2 (b).

Hyperprior feature for anchor. While introducing the position \mathbf{x}_i^k contributes to predicting the distribution of anchor features, it still lacks enough freedom to eliminate spatial redundancy. Therefore, we introduce a learnable hyperprior vector $\mathbf{z}_i \in \mathbb{R}^{50/h_c}$ for each anchor \mathbf{v}_i where h_c is a hyper-parameter to control the length of the hyperprior features. The hyperprior \mathbf{z}_i is modeled using the non-parametric, fully factorized density model [1] as follows:

$$p_{\tilde{\mathbf{z}}_i | \Theta}(\tilde{\mathbf{z}}_i | \Theta) = \prod_{j=0}^{50/h_c-1} \left(p_{z_i^j | \Theta^{(j)}}(\Theta^{(j)}) * \mathcal{U}\left(-\frac{1}{2}, \frac{1}{2}\right) \right) (\tilde{z}_i^j), \quad (9)$$

where $\tilde{\mathbf{z}}_i$ represents \mathbf{z}_i with quantization noise, j is the channel index and Θ is the network parameters for modeling the hyperprior. Since the hyperprior feature \mathbf{z}_i is quantized into integers $\hat{\mathbf{z}}$ and jointly optimized to reduce the size of the bitstream, it only occupies a small portion of storage as shown in Table 4. The final prior for coding features of the anchor \mathbf{v}_i^k is $\hat{\psi}_i^k = [\hat{\mathbf{z}}_i^k; \psi_i^k]$.

4.3 Training objective

The training objective of the proposed method is to jointly optimize the bitrate of coded anchor features and rendering loss measured by SSIM and L1 loss. The final training loss is

$$\mathcal{L} = \mathcal{L}_{scaffold} + \lambda_e \mathcal{L}_{entropy} + \lambda_m \mathcal{L}_m \quad (10)$$

where $\mathcal{L}_{scaffold}$ is the training loss of [19], \mathcal{L}_m is the masking loss from [18] to regularize the masking loss of offsets of neural Gaussians \mathbf{O}_v , and $\mathcal{L}_{entropy}$ is the overall entropy loss that measures the cost of storing anchor properties defined as follows,

$$\mathcal{L}_{entropy} = \mathbb{E}[-\log_2 p_{\tilde{\mathbf{z}}_i | \Theta}(\tilde{\mathbf{z}}_i | \Theta)] + \sum_{i=0}^{K-1} \left[-\log_2 \left[\prod_{\mathbf{f} \in \{\mathbf{f}^k, \mathbf{l}^k, \mathbf{O}^k\}} p_{\mathbf{f}}(\mathbf{f}_i | \hat{\psi}_i^k) \right] \right] \quad (11)$$

where the first term measures the cost of coding the hyperprior feature while the second term is the cost of coding features of anchor points in all the levels, and $\hat{\psi}_i^k$ is the context feature that includes both the hyperprior feature \mathbf{z}_i^k and the feature from already coded nearby anchor in level $k + 1$ as illustrated in Eq. 8.

Table 1: The quantitative results obtained from the proposed method *ContextGS* and other competitors. Baseline methods, namely 3DGS [15] and Scaffold-GS [20], are included for reference. The intermediary approaches are specifically designed for 3DGS compression. Our methodology showcases two results representing varying size and fidelity tradeoffs, achieved through adjustment of λ_e . Highlighted in red and yellow cells are the best and second-best results, respectively. Size measurements are expressed in megabytes (MB).

Datasets methods	Mip-NeRF360 [3]				Tank&Temples [16]				DeepBlending [14]				BungeeNeRF [32]			
	psnr \uparrow	ssim \uparrow	lpips \downarrow	size \downarrow	psnr \uparrow	ssim \uparrow	lpips \downarrow	size \downarrow	psnr \uparrow	ssim \uparrow	lpips \downarrow	size \downarrow	psnr \uparrow	ssim \uparrow	lpips \downarrow	size \downarrow
3DGS [15] (SIGGRAPH'23)	27.49	0.813	0.222	744.7	23.69	0.844	0.178	431.0	29.42	0.899	0.247	663.9	24.87	0.841	0.205	1616
Scaffold-GS [20] (CVPR'24)	27.50	0.806	0.252	253.9	23.96	0.853	0.177	86.50	30.21	0.906	0.254	66.00	26.62	0.865	0.241	183.0
EAGLES [11]	27.15	0.808	0.238	68.89	23.41	0.840	0.200	34.00	29.91	0.910	0.250	62.00	25.24	0.843	0.221	117.1
LightGaussian [7]	27.00	0.799	0.249	44.54	22.83	0.822	0.242	22.43	27.01	0.872	0.308	33.94	24.52	0.825	0.255	87.28
Compact3DGS [18] (CVPR'24)	27.08	0.798	0.247	48.80	23.32	0.831	0.201	39.43	29.79	0.901	0.258	43.21	23.36	0.788	0.251	82.60
Compressed3D [26] (CVPR'24)	26.98	0.801	0.238	28.80	23.32	0.832	0.194	17.28	29.38	0.898	0.253	25.30	24.13	0.802	0.245	55.79
Morgenstern <i>et al.</i> [23]	26.01	0.772	0.259	23.90	22.78	0.817	0.211	13.05	28.92	0.891	0.276	8.40	—	—	—	—
Navaneet <i>et al.</i> [25]	27.16	0.808	0.228	50.30	23.47	0.840	0.188	27.97	29.75	0.903	0.247	42.77	24.63	0.823	0.239	104.3
HAC [5]	27.53	0.807	0.238	15.26	24.04	0.846	0.187	8.10	29.98	0.902	0.269	4.35	26.48	0.845	0.250	18.49
Ours (low-rate)	27.62	0.808	0.237	12.68	24.20	0.852	0.184	7.05	30.11	0.907	0.265	3.45	26.90	0.866	0.222	14.00
Ours (high-rate)	27.75	0.811	0.231	18.41	24.29	0.855	0.176	11.80	30.39	0.909	0.258	6.60	27.15	0.875	0.205	21.80

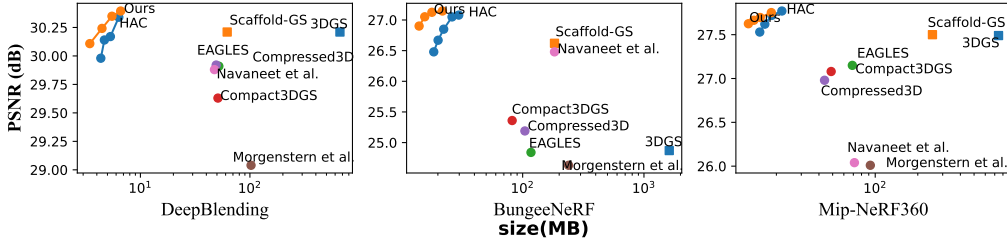


Figure 4: The Rate-Distortion (RD) curves for quantitative comparison between our method with most recent SOTA competitors. It is worth noting that the x-axis is in log scale for better visualization.

5 Experiments

5.1 Implementation details

We build our method based on Scaffold-GS [20]. The number of levels is set to 3 for all experiments and the target ratio among two adjacent iterations is 0.2. h_c is set to 4, *i.e.*, the dimension of the hyper-prior feature is a fourth of the anchor feature dimension. For a fair comparison, the dimension of the anchor feature \mathbf{f} is set to 50 following [5] and we set the same $\lambda_m = 5e - 4$. The setting of λ_e is discussed in Appendix A.3 since different values are used to evaluate different rate-distortion tradeoffs. For a fair comparison, we use the same training iterations with Scaffold-GS [20] and HAC [5], *i.e.*, 30000 iterations. Besides, we use the same hyperparameters for anchor growing as Scaffold-GS [20] so that the final model has a similar or even smaller number of anchors, leading to faster rendering speed. More implementation details are in the supplementary materials.

5.2 Comparison with baselines

Baseline, metric, and benchmark. We compare our method with 3DGS [15], Scaffold-GS [20] and some representative 3DGS compression works, including Compact3DGS [18], Compressed3D [26], EAGLES [11], LightGaussian [7], Morgenstern *et al.* [23], Navaneet *et al.* [25], and HAC [5]. The baseline methods include existing mainstream techniques, *e.g.*, pruning [7, 18], codebooks [7, 18, 25, 26], and entropy coding [11, 5, 18, 23], and includes the most recent works. We utilize PSNR, SSIM, and LPIPS [35] to evaluate the rendering qualities of different methods and report the storage size measured in MB. We evaluate the performance of the models on several real-scene datasets, including BungeeNeRF [32], DeepBlending [14], Mip-NeRF360 [3], and Tanks&Temples [16]. To more comprehensively evaluate the performance of our method, following the previous prototype [5], we use all 9 scenes in Mip-NeRF360 [3]. The detailed results of each scene are reported in the Appendix A.3. To further evaluate the performance models among a wide range of compression ratios, we use Rate-Distortion (RD) curves as an additional metric.

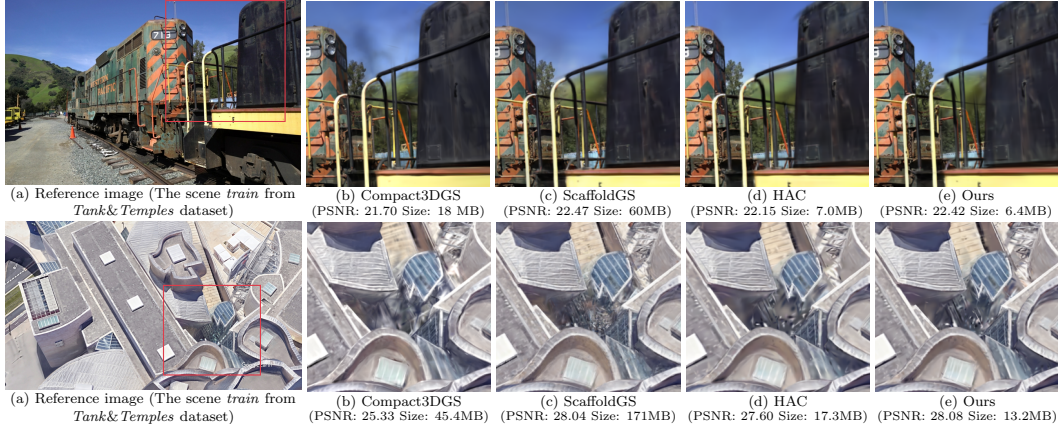


Figure 5: Visual comparisons between our method and baselines including Scaffold-GS [20], HAC [5], and Compact3DGS [18] on Bungeenerf [32] and Tank&Temples [16]. We report the PSNR (dB) of the image and the size of the 3D scene. (Best zoom in for details.)

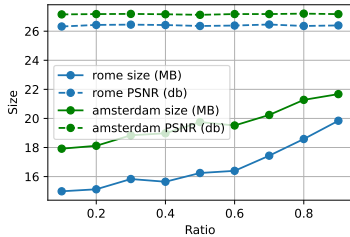


Figure 6: The ablation of different target ratio τ among different scenes. The PSNR remains relatively stable while the size of the scenes keeps increasing when increasing the τ , which demonstrates the effectiveness.

Table 2: The ablation study of each component we proposed measured on BungeeNerf [32] dataset. “HP” and “CM” represent the hyperprior and anchor level context model respectively. Ours w/o HP w/o CM can be roughly regarded as a Scaffold-GS [20] model with entropy coding and masking loss [18].

	Size (MB)	PSNR	SSIM	LPIPS
Scaffold-GS [20]	183.0	26.62	0.865	0.241
Ours w/o HP w/o CM	18.67	26.93	0.867	0.222
Ours w/o CM	15.03	26.91	0.866	0.223
Ours w/o HP	15.41	26.92	0.867	0.221
Ours	14.00	26.90	0.866	0.222

Results. As shown in Table 1, the proposed method achieves significant improvement compared to our backbone method Scaffold-GS [20] in terms of the size of the model, with a size reduction of $15\times$ in average. Besides, the proposed method also achieves higher storage efficiency compared to the most recent competitors for the 3DGS compression, *e.g.*, HAC [5], Compressed3D [26], and Compact3DGS [18]. It is worth noting that the proposed method also significantly improves rendering quality, even compared with the backbone model we use, *i.e.*, Scaffold-GS [20]. This further verifies the observation from previous works that appropriate constraints on neural Gaussians can contribute to the rendering quality, *e.g.*, entropy constraints [5] and pruning [33]. Visual comparisons between the proposed method and other competitors are shown in Fig. 5. As shown in the figure, the proposed method achieves better rendering quality with greatly reduced size compared with most recent 3D compression works and also the backbone model. Besides, a comparison of the RD curves among the proposed method and most recent competitors is shown in Fig. 4, where the proposed method achieves better performance in a wide range of compression ratios.

5.3 Ablation studies and discussions

Evaluation of target ratio τ among adjacent levels. To evaluate the performance of the proposed strategy that encodes anchors in a progressive way, we evaluate the performance of models trained under different ratios among adjacent levels. We disable the hyperprior feature to better explore the effect of different target ratios τ . As shown in Fig. 6, the PSNR remains relatively stable and the size gets relatively converged at the low ratio area. We select $\tau = 0.2$ for all experiments.

Ablation of each component. We verify the effectiveness of two main components in our methods, *i.e.*, anchor level context model and hyperprior features, and the results are shown in Table 2. We build all the models on Scaffold-GS [20] and set the model with the entropy constraint and the

Table 3: The ablation study of our method w/ and w/o reusing anchors from coarser levels, *i.e.*, anchor forward technique, measured on BungeeNerf [32] dataset.

	Size (MB)	PSNR	SSIM	LPIPS
Ours w/o dividing into levels	14.73	26.91	0.867	0.222
Ours w/o reusing anchors	15.54	26.91	0.862	0.230
Ours	13.80	26.89	0.867	0.222

Table 4: The storage cost of each component and rendering qualities of our method and baselines evaluated on the scene *rome* in BungeeNerf [32] dataset. “w/ APC” represents using anchor position coding, *i.e.*, using the hyperprior features to code the anchor positions. (The encoding/decoding time is measured on an RTX3090.)

	Number of Anchors (K)	Storage Cost (MB)								Fidelity		Speed (s)	
		Hyper	Position	Feature	Scaling	Offset	Mask	MLPs	Total	PSNR	SSIM	Encode	Decode
Scaffold-GS [20]	61.9	N/A	7.08	75.16	14.18	70.88	N/A	0.047	184.4	26.25	0.872	N/A	N/A
Ours (w/ APC)	52.3	1.026	1.954	5.708	1.603	2.556	0.452	0.320	13.62	26.38	0.871	41.33	51.58
Ours	52.5	0.778	2.543	5.808	1.586	2.563	0.452	0.316	14.06	26.38	0.871	20.40	17.85

masking loss [18] as our baseline model, *i.e.*, “Ours w/o HP w/o CM”. It is worth noting that our baseline model significantly improves the storage efficiency compared with Scaffold-GS [20] and even the latest SOTA methods. Both the proposed anchor-level context model and hyperprior feature for anchors significantly improve the compression rate compared with our strong baseline model, reducing the file size by 21% and 10.17%, respectively. Besides, using them together can further boost the performance with storage savings of 26.1% and 92.5% compared with the baseline we introduced above and Scaffold-GS [20] respectively.

Ablation of anchor forward. A main difference between the proposed method and existing works for Level-of-Detail (LOD) techniques [28] is that the proposed method can reuse the anchors of different levels, *i.e.*, anchor forward in Fig. 2 (b). For example, the anchors from different levels in [28] are stored separately. In contrast, the anchors from coarser levels are used in the final level (level 0) in our method, *i.e.*, in an autoregressive manner. To verify the effectiveness of the proposed method that reuses the anchors in coarser levels, we do an ablation study in Table 3. As shown in the table, the model w/o reusing anchors of coarser levels to the finest level leads to serious redundancy, even slightly worse than the model w/o dividing anchors into different levels. This demonstrates the effectiveness of the proposed anchor forward technique for the anchor-level context model.

Discussion on compressing anchor positions. One can utilize the hyperprior feature \mathbf{z} to predict the distribution of anchor positions and the anchor position can thereby be compressed using entropy coding. However, we find that the precision of the anchor position is essential to the performance of the model and an adaptive quantization strategy leads to serious performance degradation. While a fixed quantization width is feasible and can retain the fidelity performance while effectively compressing the size of anchors, it leads to a greatly increased number of symbols that greatly decreases the coding speed. Since the anchor position only occupies a small portion of bitstreams as shown in Table 4, we do not encode anchors into bitstreams in all the experiments.

Analysis of inference and decoding time. The rendering speed after decompression is the same as or even faster than Scaffold-GS [20] when the number of anchors is the same since we use the same data structure. However, as shown in Table 4, we can achieve higher rendering quality with fewer anchors due to the use of masking loss [18] therefore achieving faster rendering speed. For the decoding time, while the proposed method involves autoregressive coding, which is usually very slow in image compression tasks [21] due to its serial characteristics, it adds neglectable overhead to our method in both training and decompression compared to other entropy-coding-based 3DGS compression methods, such as HAC [5]. This is because, unlike autoregressive coding in image compression that predicts pixels/latent features one by one, introducing a loop of at least thousands of operations, we perform autoregressive coding group by group, introducing only a loop of 3 iterations. Additionally, there is no overlap of anchors among the coding of different levels, so the overall number of anchors to be processed remains the same as without dividing into levels.

6 Conclusion

In this work, we introduce a pioneer study into utilizing the anchor-level context model in the compression of 3D Gaussian splatting models. We divide anchors into different levels and the anchors from coarser levels are first coded and then used to predict anchors that are not coded yet. Additionally, a hyperprior feature is used for each anchor that further reduces the channel-wised redundancy. Besides, we demonstrate that utilizing the proposed anchor forward technique, *i.e.*, directly reusing the anchors from coarse levels to the final level, can achieve better performance than just using anchors of coarse levels as a prior. Extensive experiments demonstrate that the proposed methods achieve better performance than SOTA and concurrent works.

Acknowledgements

This research was done in ROSE Lab, NTU and was supported in part by the NTU-PKU Joint Research Institute, a collaboration between the Nanyang Technological University (NTU) and Peking University (PKU) that was sponsored by a donation from the Ng Teng Fong Charitable Foundation, the National Research Foundation Singapore Competitive Research Program (award number CRP29-2022-0003), and Guangdong Basic and Applied Basic Research Foundation (2024A1515010454).

References

- [1] J. Ballé, D. Minnen, S. Singh, S. J. Hwang, and N. Johnston. Variational image compression with a scale hyperprior. In *International Conference on Learning Representations*, 2018.
- [2] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021.
- [3] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.
- [4] A. Chen, Z. Xu, A. Geiger, J. Yu, and H. Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, pages 333–350. Springer, 2022.
- [5] Y. Chen, Q. Wu, J. Cai, M. Harandi, and W. Lin. Hac: Hash-grid assisted context for 3d gaussian splatting compression. *arXiv preprint arXiv:2403.14530*, 2024.
- [6] K. Deng, A. Liu, J.-Y. Zhu, and D. Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022.
- [7] Z. Fan, K. Wang, K. Wen, Z. Zhu, D. Xu, and Z. Wang. Lightgaussian: Unbounded 3d gaussian compression with 15x reduction and 200+ fps. *arXiv preprint arXiv:2311.17245*, 2023.
- [8] S. Fridovich-Keil, G. Meanti, F. R. Warburg, B. Recht, and A. Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023.
- [9] S. Fridovich-Keil, A. Yu, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022.
- [10] Q. Gao, Q. Xu, H. Su, U. Neumann, and Z. Xu. Strivec: Sparse tri-vector radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17569–17579, 2023.
- [11] S. Girish, K. Gupta, and A. Shrivastava. Eagles: Efficient accelerated 3d gaussians with lightweight encodings. *arXiv preprint arXiv:2312.04564*, 2023.
- [12] S. Girish, A. Shrivastava, and K. Gupta. Shacira: Scalable hash-grid compression for implicit neural representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17513–17524, 2023.
- [13] K. Han and W. Xiang. Multiscale tensor decomposition and rendering equation encoding for view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4232–4241, 2023.

- [14] P. Hedman, J. Philip, T. Price, J.-M. Frahm, G. Drettakis, and G. Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (ToG)*, 37(6):1–15, 2018.
- [15] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023.
- [16] A. Knapitsch, J. Park, Q.-Y. Zhou, and V. Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017.
- [17] C. Lassner and M. Zollhofer. Pulsar: Efficient sphere-based neural rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1440–1449, 2021.
- [18] J. C. Lee, D. Rho, X. Sun, J. H. Ko, and E. Park. Compact 3d gaussian representation for radiance field. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [19] T. Lu, M. Yu, L. Xu, Y. Xiangli, L. Wang, D. Lin, and B. Dai. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. *arXiv preprint arXiv:2312.00109*, 2023.
- [20] T. Lu, M. Yu, L. Xu, Y. Xiangli, L. Wang, D. Lin, and B. Dai. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [21] F. Mentzer, E. Agustsson, M. Tschannen, R. Timofte, and L. Van Gool. Conditional probability models for deep image compression. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4394–4402, 2018.
- [22] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- [23] W. Morgenstern, F. Barthel, A. Hilsmann, and P. Eisert. Compact 3d scene representation via self-organizing gaussian grids. *arXiv preprint arXiv:2312.13299*, 2023.
- [24] T. Müller, A. Evans, C. Schied, and A. Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022.
- [25] K. Navaneet, K. P. Meibodi, S. A. Koochpayegani, and H. Pirsiavash. Compact3d: Compressing gaussian splat radiance field models with vector quantization. *arXiv preprint arXiv:2311.18159*, 2023.
- [26] S. Niedermayr, J. Stumpfegger, and R. Westermann. Compressed 3d gaussian splatting for accelerated novel view synthesis. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- [27] C. Reiser, S. Peng, Y. Liao, and A. Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 14335–14345, 2021.
- [28] K. Ren, L. Jiang, T. Lu, M. Yu, L. Xu, Z. Ni, and B. Dai. Octree-gs: Towards consistent real-time rendering with lod-structured 3d gaussians. *arXiv preprint arXiv:2403.17898*, 2024.
- [29] D. Rho, B. Lee, S. Nam, J. C. Lee, J. H. Ko, and E. Park. Masked wavelet representation for compact neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20680–20690, 2023.
- [30] T. Takikawa, A. Evans, J. Tremblay, T. Müller, M. McGuire, A. Jacobson, and S. Fidler. Variable bitrate neural fields. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–9, 2022.
- [31] A. Van den Oord, N. Kalchbrenner, L. Espeholt, O. Vinyals, A. Graves, et al. Conditional image generation with pixelcnn decoders. *Advances in neural information processing systems*, 29, 2016.
- [32] Y. Xiangli, L. Xu, X. Pan, N. Zhao, A. Rao, C. Theobalt, B. Dai, and D. Lin. Bungeenerf: Progressive neural radiance field for extreme multi-scale scene rendering. In *European conference on computer vision*, pages 106–122. Springer, 2022.
- [33] R. Yang, Z. Zhu, Z. Jiang, B. Ye, X. Chen, Y. Zhang, Y. Chen, J. Zhao, and H. Zhao. Spectrally pruned gaussian fields with neural compensation. *arXiv preprint arXiv:2405.00676*, 2024.
- [34] Q. Zhang, S.-H. Baek, S. Rusinkiewicz, and F. Heide. Differentiable point-based radiance fields for efficient view synthesis. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–12, 2022.
- [35] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.

A Appendix / supplemental material

A.1 Broader impacts

The social impacts of our work are mainly in three folds:

1. Accessibility: This work contributes to making advanced 3DGS [15, 20] models more accessible to a wider audience. The reduced size of the 3DGS representation means that it can be more easily stored, transmitted, and processed on various devices, potentially democratizing access to high-quality rendering capabilities.
2. Applications: Improved compression techniques for 3DGS could enhance various applications across industries, including virtual reality, gaming, medical imaging, and architectural visualization. These advancements may lead to more immersive experiences, better medical diagnostics, and more efficient design workflows.
3. Open Access: The commitment to releasing the code fosters transparency and collaboration within the research community. Open access to the code allows other researchers and practitioners to build upon this work, accelerating innovation in the field of view synthesis and compression.

We do not find serious negative impacts since our work is only for compression. There is no concerns regarding the misuse of our models or generating fake data.

A.2 Limitations

A main and inevitable limitation of the proposed method is that the entropy coding process introduces extra computational costs to estimate the entropy of the anchor features during training encoding/decoding when saving/loading the 3D scene. For example, it requires extra time to decode the data of 3D scenes from the bitstream, making it challenging to start rendering at once when clicking the file.

A.3 More experimental details and results

We report the detailed comparison of our method w/ and w/o using the coding of anchor position in Table 5. To evaluate the performance among different rate-distortion (RD) tradeoffs, we utilize different λ_e , *i.e.*, a larger λ_e leads to a smaller model size. For a fair comparison, we use the same normalization of the hyper-parameter λ_e as [5] by additionally utilizing the dimension of anchor features as the divisor so that the same λ_e can lead to the similar bias towards the RD tradeoffs. The detailed results of each scene with different λ_e are reported in Table 6, Table 8, Table 7, and Table 9.

Table 5: Quantitative results of the effect of anchor position coding on our method. For our approach, we give two results of different size and fidelity tradeoffs by adjusting λ_e . A smaller λ_e results in a larger size but improved fidelity, and vice versa.

Datasets methods	Mip-NeRF360 [3]				Tank&Temples [16]				DeepBlending [14]				BungeeNeRF [32]			
	psnr↑	ssim↑	lpips↓	size↓	psnr↑	ssim↑	lpips↓	size↓	psnr↑	ssim↑	lpips↓	size↓	psnr↑	ssim↑	lpips↓	size↓
w/ encoding anchors																
Ours (low-rate)	27.61	0.809	0.236	11.32	24.16	0.851	0.185	6.91	30.11	0.907	0.269	3.31	26.89	0.867	0.222	13.80
Ours (high-rate)	27.86	0.813	0.230	21.07	24.29	0.855	0.178	11.47	30.42	0.910	0.261	6.40	27.15	0.876	0.202	25.23
w/o encoding anchors																
Ours (low-rate)	27.62	0.808	0.237	12.68	24.20	0.852	0.184	7.05	30.11	0.907	0.265	3.45	26.90	0.866	0.222	14.00
Ours (high-rate)	27.75	0.811	0.231	18.41	24.29	0.855	0.176	11.80	30.39	0.909	0.258	6.60	27.15	0.875	0.205	21.80

Table 6: Per-scene results of our method on BungeeNerf [32] dataset.

	Scene	Size	PSNR	SSIM	LPIPS
low-rate ($\lambda = 0.004$)	rome	13.9631	25.9751	0.8669	0.2143
	quebec	11.5909	30.1941	0.9337	0.1669
	pompidou	15.2258	25.4676	0.8469	0.2446
	hollywood	13.5183	24.5154	0.7772	0.3164
	bilbao	13.1537	28.0842	0.8877	0.1932
	amsterdam	16.5549	27.1694	0.8842	0.1956
	Average	14.0011	26.9010	0.8661	0.2218
high-rate ($\lambda = 0.001$)	rome	21.7033	26.6297	0.8825	0.1929
	quebec	18.1517	30.5271	0.9400	0.1510
	pompidou	23.4748	25.6218	0.8546	0.2330
	hollywood	20.7082	24.7170	0.7876	0.3002
	bilbao	20.7686	27.9842	0.8928	0.1770
	amsterdam	26.0079	27.3979	0.8949	0.1754
	Average	21.8024	27.1463	0.8754	0.2049

Table 7: Per-scene results of our method on DeepBlending [14] dataset.

	Scene	Size	PSNR	SSIM	LPIPS
low-rate ($\lambda = 0.004$)	drjohnson	3.94	29.68	0.906	0.261
	playroom	2.96	30.53	0.907	0.269
	Average	3.45	30.11	0.907	0.265
high-rate ($\lambda = 0.0005$)	drjohnson	7.80	29.86	0.909	0.249
	playroom	5.41	30.93	0.910	0.268
	Average	6.60	30.39	0.909	0.258

Table 8: Per-scene results of our method on Tank&Template [16] dataset.

	Scene	Size	PSNR	SSIM	LPIPS
low-rate ($\lambda = 0.004$)	train	6.39	22.40	0.818	0.217
	truck	7.72	26.00	0.885	0.150
	Average	7.05	24.20	0.852	0.184
high-rate ($\lambda = 0.0005$)	train	10.55	22.53	0.823	0.208
	truck	13.06	26.05	0.888	0.143
	Average	11.80	24.29	0.855	0.176

Table 9: Per-scene results of our method on the Mip-NeRF360 [3] dataset.

	Scene	Size	PSNR	SSIM	LPIPS
low-rate ($\lambda = 0.004$)	bicycle	21.82	25.08	0.738	0.271
	bonsai	7.17	32.67	0.946	0.186
	counter	6.30	29.38	0.912	0.197
	flowers	16.71	21.31	0.576	0.377
	garden	18.78	27.32	0.845	0.148
	kitchen	7.00	31.28	0.925	0.131
	room	4.50	31.71	0.923	0.207
	stump	14.86	26.58	0.762	0.268
	treehill	17.00	23.29	0.647	0.349
	Average	12.68	27.62	0.808	0.237
high-rate ($\lambda = 0.0005$)	bicycle	38.09	24.97	0.740	0.263
	bonsai	12.43	32.93	0.951	0.182
	counter	10.67	29.60	0.916	0.190
	flowers	28.27	21.18	0.571	0.378
	garden	31.26	27.39	0.851	0.136
	kitchen	12.44	31.69	0.931	0.123
	room	8.09	32.03	0.928	0.196
	stump	24.22	26.56	0.761	0.263
	treehill	28.78	23.09	0.645	0.346
	Average	21.58	27.72	0.811	0.231

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: We clearly summarized the contributions in both abstract and introduction.

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: We discuss the limitations in the Appendix A.2

Guidelines:

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

3. Theory Assumptions and Proofs

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

Answer: [NA]

Justification: We do 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: We included all necessary information to reproduce the main results. Besides, we will release the code upon acceptance.

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

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We will release the code and pretrained models on acceptance. The data we used are all public data.

Guidelines:

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

6. Experimental Setting/Details

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

Answer: [Yes]

Justification: The code will be released on acceptance. We included the necessary details for training and testing.

Guidelines:

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

7. Experiment Statistical Significance

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

Answer: [No]

Justification: We evaluate our methods over different datasets and rate-distortion tradeoffs, resulting in a large set of experiments. The proposed method outperform other SOTA methods by a large margin, which is far beyond the standard deviation of the experiments.

Guidelines:

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

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

8. Experiments Compute Resources

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

Answer: [Yes]

Justification: All our experiments are done using a server with RTX3090s. The requirements of compute resources are small.

Guidelines:

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

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: Our work does not have special ethics concerns. All data and experiments comply with ethical guidelines.

Guidelines:

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

10. Broader Impacts

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

Answer: [Yes]

Justification: We discuss broader impacts in Appendix A.1

Guidelines:

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

- 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).

11. Safeguards

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

Answer: [NA]

Justification: Our method does not have a high risk of misuse.

Guidelines:

- 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.

12. Licenses for existing assets

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

Answer: [Yes]

Justification: We cite the original papers and state which version of the dataset we used.

Guidelines:

- 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.
- 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.

13. **New Assets**

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

Answer: [NA]

Justification: We do not introduce new data.

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.

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

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

Answer: [NA]

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

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.

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

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

Answer: [NA]

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

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.
- 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.