A Hierarchical Compression Technique for 3D Gaussian Splatting Compression

He Huang* Wenjie Huang* Qi Yang[†] Yiling Xu* Zhu Li[†]

* Shanghai Jiao Tong University, † University of Missouri-Kansas City {huanghe0429, huangwenjie2023, yl.xu}@sjtu.edu.cn, littlleempty@gmail.com, lizhu@umkc.edu

Abstract—3D Gaussian Splatting (GS) demonstrates excellent rendering quality and generation speed in novel view synthesis. However, substantial data size poses challenges for storage and transmission, making 3D GS compression an essential technology. Current 3D GS compression research primarily focuses on developing more compact scene representations, such as converting explicit 3D GS data into implicit forms. In contrast, compression of the GS data itself has hardly been explored. To address this gap, we propose a Hierarchical GS Compression (HGSC) technique. Initially, we prune unimportant Gaussians based on importance scores derived from both global and local significance, effectively reducing redundancy while maintaining visual quality. An Octree structure is used to compress 3D positions. Based on the 3D GS Octree, we implement a hierarchical attribute compression strategy by employing a KD-tree to partition the 3D GS into multiple blocks. We apply farthest point sampling to select anchor primitives within each block and others as non-anchor primitives with varying Levels of Details (LoDs). Anchor primitives serve as reference points for predicting nonanchor primitives across different LoDs to reduce spatial redundancy. For anchor primitives, we use the region adaptive hierarchical transform to achieve near-lossless compression of various attributes. For non-anchor primitives, each is predicted based on the k-nearest anchor primitives. To further minimize prediction errors, the reconstructed LoD and anchor primitives are combined to form new anchor primitives to predict the next LoD. Our method notably achieves superior compression quality and a significant data size reduction of over 4.5× compared to the state-of-the-art compression method on small scenes datasets.

Index Terms-3D Gaussian Splatting, Compression

I. INTRODUCTION

3D Gaussian Splatting (GS) [1] has demonstrated substantial advances in the field of novel view synthesis due to its impressive visual quality with ultra fast training speed. Different from Neural Radiance Field (NeRF) [2] with implicit representations, 3D GS uses serial explicit scattered isotropic ellipsoids to reconstruct the 3D scene. Each Gaussian consists of a 3D point center and several attributes, including scale vector, rotation quaternion, Spherical Harmonic (SH) coefficients, and opacity. Leveraging a highly optimized CUDA-based rendering implementation, 3D GS enables rapid training and rendering, making it highly suitable for practical applications. Additionally, its explicit data format is not only easy to understand and analyze, but also facilitates downstream processing (e.g., MGA [3]), making it a promising candidate for industrial application and standardization efforts.

However, explicit point-based representations inherently result in significant storage overhead, as each point and its associated attributes must be stored independently. For example, reconstructing a large scene typically requires several million Gaussians, which can consume more than one gigabyte of memory. Therefore, compression of 3D GS becomes an essential technology to mitigate storage and transmission overhead.

Currently, research on 3D GS compression can be categorized into two distinct branches: generative compression and traditional compression [4]. The majority of studies focus on generative compression, employing techniques such as pruning, codebooks, and entropy constraints to produce more compact data representations during 3D GS generation. For example, LightGaussian [5] prunes insignificant Gaussians based solely on the opacity parameter, and HAC [6] introduces a hash-grid assisted context model to reduce spatial redundancy. However, although traditional compression has hardly been studied now, it remains an equally important area of research. GGSC [4] is the first work to address this gap by proposing a simple but effective graph-based compression anchor. However, the performance of GGSC is limited as it does not fully exploit the spatial redundancy of 3D GS.

In this paper, we introduce a Hierarchical GS Compression (HGSC) technique. We first prune Gaussians based on primitive importance scores assessed by both global significance and local significance: global significance is determined by each Gaussian's contribution to rendering color across different views, while local significance is associated with the volume of each Gaussian, both of which are crucial for maintaining final rendering quality. Then, an Octree [7] structure is employed to compress 3D positions. To reduce the influence of point merging within a voxel in Octree, the reconstructed points are recolored by applying the attributes of the nearest Gaussian from original 3D GS to ensure consistency. Based on the 3D GS Octree, we implement a hierarchical compression strategy. Specifically, we use a KD-tree [8] to split the 3D GS into multiple blocks and apply Farthest Point Sampling (FPS) [9] to select anchor primitives in each block and then generate varying Levels of Details (LoDs) primitives. These anchor primitives serve as references for predicting non-anchor primitives across different LoDs to reduce spatial redundancy. For anchor primitives, we employ the region adaptive hierarchical transform (RAHT) [10] to achieve near-lossless compression of various attributes to enhance prediction ac-

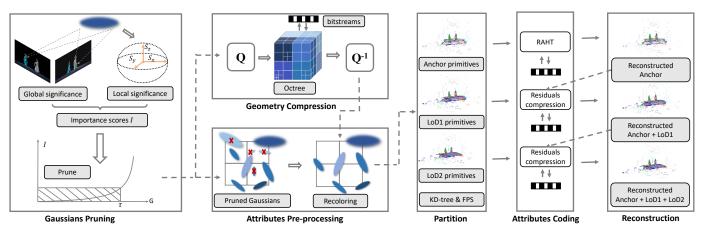


Fig. 1: Framework of HGSC. Q and Q^{-1} denote the processes of quantization and dequantization, respectively.

curacy. For non-anchor primitives, each is predicted by the k-nearest anchor primitives. Subsequently, the discrepancies between the predicted and actual attributes are quantized and subsequently encoded using the LZ77 [11] codec. To minimize prediction errors, the current reconstructed LoD and anchor primitives are combined to form the new anchor primitives for predicting the next LoD. Overall, our method achieves better compression efficiency and reduced processing time compared to the benchmark GGSC.

II. PROPOSED METHODS

A. Preliminary

3D GS is an explicit point-based 3D scene representation through a differentiable splatting and tile-based rasterization. Each Gaussian geometry is characterized by a point center X and covariance matrix Σ ,

$$G(X) = e^{-\frac{1}{2}X^T \Sigma^{-1} X}.$$
 (1)

To maintain the positive semi-definite characteristics, Σ is decomposed into a scaling matrix S = diag(s) and a rotation matrix $R: \Sigma = RSS^TR^T$. To render an image from a random viewpoint, 3D Gaussians are first splatted to 2D planes, and the pixel value $C \in \mathbb{R}^3$ is computed using α -blending,

$$C = \sum_{i \in N} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j),$$
 (2)

where α measures the opacity of each Gaussian after 2D projection, $c \in \mathbb{R}^3$ is view-dependent color modeled by Spherical Harmonic (SH) coefficients, and N is the number of sorted Gaussians contributing to the rendering.

B. overview

We propose an efficient hierarchical 3D GS compression technique. As depicted in Fig. 1, unimportant Gaussians are initially pruned based on assessments of global and local significance. For geometry, an Octree structure is employed to compress the 3D positions. For attributes, a hierarchical compression strategy is proposed. Specifically, the anchor primitives and different LoDs are first partitioned using KD-tree and FPS techniques. Given that anchor primitives serve as references, we subsequently apply RAHT for near-lossless

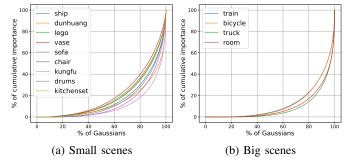


Fig. 2: Cumulative distribution curves of importance. compression to ensure fundamental information and improve prediction accuracy. For each LoD, the attributes of each non-anchor primitives are predicted based on the k-nearest anchor primitives and lower LoDs. Subsequently, the discrepancies between the predicted and actual attributes are quantized and subsequently encoded. After completion of the compression of the current LoD, the reconstructed LoD and anchor primitives are merged to form the new anchor primitives for the subsequent LoDs, enhancing prediction accuracy. This iterative process continues until the entire 3D GS is effectively compressed.

C. Gaussians Pruning

3D GS contains a large number of Gaussians and many of them are relatively insignificant for visual quality. To improve compression efficiency, we prune these less important Gaussians based on importance scores that involve two components: global significance I_g and local significance I_l . The overall importance score of a Gaussian is defined as $I = I_g I_l$, with both components being crucial for achieving high-quality final rendering. Based on (2), we define global significance I_g as:

$$I_g = \sum_{p \in \mathcal{P}} \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \tag{3}$$

where \mathcal{P} represents the set of pixels that are overlapped by the projection of the Gaussian, while i denotes the set of sorted Gaussians along the ray which is determined by the α -blending. The local significance I_l is defined as:

$$I_l = (V_{norm})^{\beta}, V_{norm} = \min\left(\max\left(\frac{V_G}{V_{max90}}, 0\right), 1\right).$$
 (4)

Here, the Gaussian volume, which is the product of the scale vector: $V_G = S_x S_y S_z$, is firstly normalized by the 90% largest of all sorted Gaussians. Then, it is clipped to the range between 0 and 1 to prevent excessive floating Gaussians. The β is introduced to provide additional flexibility.

As shown in Fig. 2, we sort the importance scores and visualize their cumulative distribution curves. Approximately 60% of the Gaussians only account for less than 20% of the total importance in both small and large scenes, where the importance refers to the contribution to the final rendering results. Hence, we prune the $\tau\%$ of Gaussians with the lowest importance scores.

D. Geometry Compression

After pruning, we compress the 3D positions with Octree. The whole space is recursively divided into 8 subvoxels up to the depth of *d*. The occupancy symbol for each voxel is composed of 8 bits (1 to 255 in decimal), where each bit indicates the occupancy status of the corresponding subvoxel. Subsequently, the mutual information among subvoxels is utilized for advanced context modeling, and each occupancy symbol is encoded using an arithmetic encoder [12].

E. Attributes Compression

Attributes pre-processing. After geometry compression, multiple points within a voxel in Octree are merged into a single point. To maintain color consistency, the reconstructed points are recolored by applying the attributes of the nearest Gaussian from the original 3D GS. Subsequently, inspired by image compression techniques, we convert the Spherical Harmonic (SH) coefficients from the RGB to the YUV color space. This conversion is based on the principle that the Luminance Component (Y) is more critical for visual fidelity than the Chrominance components (U and V), allowing for more flexible compression parameter settings.

Anchor primitives and different LoDs Partition. Based on a comprehensive visual analysis, we have identified spatial redundancy among attributes within the 3D GS. To address this, we propose sampling anchor primitives to facilitate predictions for neighboring non-anchor primitives. However, this method can be prone to significant prediction errors. Therefore, we further introduce a hierarchical compression strategy. Initially, the anchor primitives are sampled across the entire 3D GS to serve as reference points. Non-anchor primitives in higher LoDs are then predicted based on the anchor primitives and lower LoDs, which are generated using the same approach as for the anchor primitives.

Considering Octree may create numerous empty cubes due to the uneven spatial distribution of Gaussians, we employ KD-tree to divide the 3D GS into blocks, ensuring a more uniform distribution of points within each block. To ensure that the anchor primitives provide comprehensive coverage as reference points, we apply FPS within each block to select the anchor primitives. Subsequently, different LoDs are generated from the remaining non-anchor primitives by FPS across the blocks.

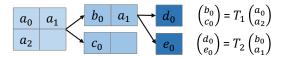


Fig. 3: 2D example of RAHT

RATH for anchor primitives. We posit that anchor primitives are of critical importance, as they serve as reference for predicting non-anchor primitives. To preserve accuracy, RAHT is used for near-lossless compression. RAHT predicts the attributes of voxels at a higher level of the Octree by using the attributes associated with voxels from the lower level. As shown in Fig. 3, RAHT merge voxels along vertical direction with transform T_1 ,

$$T_1 = \frac{1}{\sqrt{w_1} + \sqrt{w_2}} \begin{bmatrix} \sqrt{w_1} & \sqrt{w_2} \\ -\sqrt{w_2} & \sqrt{w_1} \end{bmatrix},$$
 (5)

where the weight coefficient w_i is the number of primitives that corresponding voxel contains. Similarly to T_1 , the coefficients b_0 and a_1 are merged with the transform T_2 . Ultimately, one direct current (DC) coefficient d_0 and two alternating current (AC) coefficients c_0 and e_0 are acquired and encoded by arithmetic encoding [12]. The decoding process is the inverse transform.

Residuals compression for LoDs primitives. For each LoD, each primitive is predicted by the k-nearest anchor primitives. The residuals between the predicted attributes and the true attributes are quantized and then encoded by the LZ77 codec. To reduce predicted errors, the current reconstructed LoD and anchor primitives are combined as the new anchor primitives to predict the next LoD.

III. EXPERIMENT RESULTS

A. Experimental Setting

Consistent with GGSC, we evaluate our method on two types of datasets. The small scenes often need less than one million Gaussians, including 11 different datasets from [2], [13], [14]. Meanwhile, the big scenes require millions of Gaussians, including 4 different datasets from [15]–[17]. We evaluate rendering quality and compression performance using PSNR, SSIM [18], LPIPS [19], size, encoding and decoding time compared to benchmark. For the hyperparameter settings, we set the depth of the Octree to 12. The pruning threshold $\tau\%$ is set to 60% for small scenes dataset and 66% for big scenes dataset. Additionally, the number of LoDs is set to 2, with 10% of the total primitives sampled as anchor primitives, 30% for LoD1 and 60% for LoD2.

B. Results

The quantitative results in Tab. I and the qualitative outputs in Fig. 5 clearly demonstrate the effectiveness of our proposed HGSC method. HGSC achieves a substantial size reduction of over $4.5 \times \text{compared}$ to the vanilla 3D GS, without any noticeable loss in rendering quality. We provide two results of different bitrate by adjusting τ and quantization bit depths, which consistently achieve the best or second-best results across most metrics on all datasets. To further ensure a fair

TABLE I: Quantitative results of our approach and benchmark on small and big scenes datasets. 3D GS [1] is the baseline. For our approach, we give two results of different size and fidelity trade-offs by adjusting τ and quantization bit depths. A smaller τ and larger bit depths results in a larger size but improved fidelity, and vice versa. The best and 2nd best results are highlighted in red and yellow cells. E-time and D-time are short for encoding and decoding time respectively.

Methods	Small Scenes Datasets					Big Scenes Datasets						
	Size↓ (MB)	PSNR↑	SSIM↑	LPIPS↓	E-time↓ (s)	D-time↓ (s)	Size↓ (MB)	PSNR↑	SSIM↑	LPIPS↓	E-time↓ (s)	D-time↓ (s)
3D GS [1]	49.74	-	-	-	=	-	253.33	-	-	-	-	-
GGSC-lowrate	16.10	29.20	0.923	0.072	119.83	65.14	42.98	19.06	0.672	0.357	381.74	45.35
GGSC-highrate	20.39	39.95	0.983	0.024	383.81	166.08	78.43	24.24	0.878	0.167	1212.88	117.63
Ours-lowrate	11.45	38.23	0.989	0.013	5.23	1.95	41.70	24.64	0.864	0.158	19.92	9.06
Ours-highrate	15.81	41.31	0.994	0.007	6.37	2.14	66.85	26.00	0.897	0.118	49.53	23.97

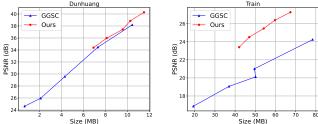


Fig. 4: RD curves for quantitative comparisons on "Dunhuang" from [14] and "Train" from [17].

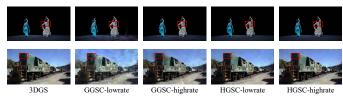


Fig. 5: Qualitative comparisons of "dunhuang" (Top) and "train" (Bottom).

comparison, we present Rate-Distortion (RD) curves in Fig. 4, showing that our method achieves at least a 6% BD-rate [20] improvement. This improvement can be attributed to two key factors: 1) the precise calculation of importance scores, which allows for the pruning of insignificant primitives, and 2) our hierarchical compression strategy, which effectively minimizes spatial redundancy. While GGSC manages some size reduction, it introduces substantial distortions in rendering results and leads to a considerable drop in objective quality. Moreover, our method boasts rapid decoding times, completing the process in approximately 2 seconds.

IV. ABLATION STUDY

As shown in Tab. II, we conducted ablation studies to evaluate the impact of each component in our method. The results demonstrate that Gaussians Pruning effectively removes redundant Gaussians, retaining only those that are essential for accurate scene representation. The Geometry Compression module offers a modest reduction in size, as geometry bit-streams constitute a minor portion of the total data. Significant size reduction is achieved through Attributes Compression by employing a hierarchical compression strategy that effectively reduces spatial redundancy in the attribute data. Furthermore, the SH Transformation module improves rendering quality by transforming RGB data into the YUV domain, which is more conducive to compression. Collectively, these modules

TABLE II: Ablation study of different stages on "Dunhuang". "+" indicates adding current model to the previous stage.

Stages	Size↓ (MB)	PSNR↑	SSIM↑	LPIPS↓
Baseline (3D GS [1])	37.75	-	-	-
+ Gaussians Pruning	14.37	41.20	0.9952	0.0063
+ Geometry Compression	13.58	41.08	0.9946	0.0075
+ Attributes Compression	8.12	39.78	0.9926	0.0117
+ SH Transformation	8.12	40.14	0.9937	0.0105

TABLE III: Comparison of pruning strategy on "Train".

Methods	PSNR↑	SSIM↑	LPIPS↓
LightGaussian [5]	20.50	0.8417	0.1625
Ours (I_g)	26.52	0.8976	0.1215
Ours (I_l)	16.79	0.7452	0.2254
Ours (I_gI_l)	27.26	0.9111	0.1053

substantially enhance the overall efficiency and performance of our framework. In addition, as illustrated in Tab. III, our proposed pruning strategy outperforms LightGaussian [5], which prunes Gaussians based solely on the opacity parameter. The results emphasize the importance of incorporating both global significance I_g and local significance I_l in the pruning process to maximize the retention of critical information. Notably, the output sizes of these strategies are consistent because they prune the same percentage of Gaussians, thereby maintaining a balanced approach to data reduction on 15 different datasets.

V. CONCLUSION

In this paper, we propose a Hierarchical 3D Gaussian Splatting Compression (HGSC) technique to enhance traditional GS compression. Initially, we prune unimportant Gaussian primitives based on both global and local significance. Then, an Octree structure is used to compress 3D positions. Following the 3D GS Octree, we implement a hierarchical compression strategy for attributes. Specifically, we use a KD-tree to split the 3DGS into multiple blocks and apply FPS to select anchor primitives in each block and non-anchor primitives within varying Levels of Details (LoDs). For anchor primitives, we employ RAHT to achieve near-lossless compression of various attributes. For non-anchor primitives, each is predicted by the k-nearest anchor. To minimize prediction errors, the current reconstructed LoD and anchor primitives are combined to form the new anchor primitives to predict the next LoD. Our method has achieved new SOTA compression performance and reconstruction rendering quality over concurrent works on 15 different datasets. Future work will focus on further reducing the time complexity.

REFERENCES

- B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering." ACM Trans. Graph., vol. 42, no. 4, pp. 139–1, 2023.
- [2] 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*, vol. 65, no. 1, pp. 99–106, 2021.
- [3] Y.-C. Sun, Y. Shi, W. T. Ooi, C.-Y. Huang, and C.-H. Hsu, "Multi-frame bitrate allocation of dynamic 3d gaussian splatting streaming over dynamic networks," in *Proceedings of the 2024 SIGCOMM Workshop on Emerging Multimedia Systems*, 2024, pp. 1–7.
- [4] Q. Yang, K. Yang, Y. Xing, Y. Xu, and Z. Li, "A benchmark for gaussian splatting compression and quality assessment study," arXiv preprint arXiv:2407.14197, 2024.
- [5] 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.
- [6] 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.
- [7] D. Meagher, "Octree encoding: A new technique for the representation, manipulation and display of arbitrary 3-d objects by computer," *Rens-selaer Polytechnic Institute*, 1980.
- [8] J. L. Bentley, "Multidimensional binary search trees used for associative searching," *Communications of the ACM*, vol. 18, no. 9, pp. 509–517, 1075
- [9] P. S. Heckbert and M. Garland, "The farthest point strategy for progressive mesh simplification," in *Proceedings of the 24th annual conference on Computer graphics and interactive techniques*. ACM Press/Addison-Wesley Publishing Co., 1997, pp. 209–216.
- [10] R. L. De Queiroz and P. A. Chou, "Compression of 3d point clouds using a region-adaptive hierarchical transform," *IEEE Transactions on Image Processing*, vol. 25, no. 8, pp. 3947–3956, 2016.
- [11] J. Ziv and A. Lempel, "Compression of individual sequences via variable-rate coding," *IEEE transactions on Information Theory*, vol. 24, no. 5, pp. 530–536, 1978.
- [12] J. J. Rissanen, "Generalized kraft inequality and arithmetic coding," *IBM Journal of research and development*, vol. 20, no. 3, pp. 198–203, 1976.
- [13] Y. Xing, Q. Yang, K. Yang, Y. Xu, and Z. Li, "Explicit_nerf_qa: A quality assessment database for explicit nerf model compression," arXiv preprint arXiv:2407.08165, 2024.
- [14] X. Zheng, L. Liao, X. Li, J. Jiao, R. Wang, F. Gao, S. Wang, and R. Wang, "Pku-dymvhumans: A multi-view video benchmark for highfidelity dynamic human modeling," in *Proceedings of the IEEE/CVF* Conference on Computer Vision and Pattern Recognition, 2024, pp. 22 530–22 540.
- [15] 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*, 2022, pp. 5470–5479.
- [16] 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), vol. 37, no. 6, pp. 1–15, 2018.
- [17] A. Knapitsch, J. Park, Q.-Y. Zhou, and V. Koltun, "Tanks and temples: Benchmarking large-scale scene reconstruction," ACM Transactions on Graphics (ToG), vol. 36, no. 4, pp. 1–13, 2017.
- [18] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [19] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 586–595.
- [20] G. Bjontegaard, "Calculation of average psnr differences between rdcurves," ITU SG16 Doc. VCEG-M33, 2001.