# FlexGaussian: Flexible and Cost-Effective Training-Free Compression for 3D Gaussian Splatting

Boyuan Tian
boyuant2@illinois.edu
UIUC

Qizhe Gao
qizheg2@illinois.edu
UIUC

Siran Xianyu
sxianyu2@illinois.edu
UIUC

Xiaotong Cui
xcui15@illinois.edu
UIUC

Minjia Zhang
minjiaz@illinois.edu
UIUC

## Abstract

3D Gaussian splatting has become a prominent technique for representing and rendering complex 3D scenes, due to its high fidelity and speed advantages. However, the growing demand for large-scale models calls for effective compression to reduce memory and computation costs, especially on mobile and edge devices with limited resources. Existing compression methods effectively reduce 3D Gaussian parameters but often require extensive retraining or fine-tuning, lacking flexibility under varying compression constraints.

In this paper, we introduce FlexGaussian, a flexible and cost-effective method that combines mixed-precision quantization with attribute-discriminative pruning for training-free 3D Gaussian compression. FlexGaussian eliminates the need for retraining and adapts easily to diverse compression targets. Evaluation results show that FlexGaussian achieves up to 96.4% compression while maintaining high rendering quality (<1 dB drop in PSNR), and is deployable on mobile devices. FlexGaussian delivers high compression ratios within seconds, being 1.7-2.1× faster than state-of-the-art training-free methods and 10-100× faster than training-involved approaches. The code is being prepared and will be released soon at: https://github.com/Supercomputing-System-AI-Lab/FlexGaussian

## Keywords

3DGS, Training-free Compression, Pruning, Quantization, Mobile

## 1 Introduction

3D Gaussian splatting (3D-GS) has emerged as a powerful technique for novel view synthesis (NVS) [13, 23], which generates new views of a 3D scene by interpolating from a small set of images with known camera parameters. By leveraging point-based representations (e.g., Gaussian primitives) to model scene geometry and appearance, 3D-GS delivers enhanced realism, detailed lighting information, and significantly faster rendering speeds than alternative Neural Radiance Fields (NeRF) methods [2, 23, 24], making it a promising candidate for enhancing applications in mobile, gaming, and AR/VR devices. However, the heavy storage requirements for Gaussians often hinder the deployment of 3D-GS models on hardware platforms with limited resources, particularly when representing large-scale, high-resolution scenes that naturally involve more Gaussians. For intuitive illustration, the GARDEN scene from the Mip-NeRF360 dataset, covering an area of less than $100m^2$ and trained at $1.6K$ image resolution, requires over 5.8 million explicit Gaussian primitives, demanding gigabyte level of storage [13].

To facilitate the deployment of 3D-GS, prior works proposed compression pipelines to minimize the number of Gaussians [7, 10, 27, 29, 37] or create compact representations for 3D-GS [5, 22, 26] while preserving rendering quality. Despite promising results, existing methods either modify the training procedure for compactness
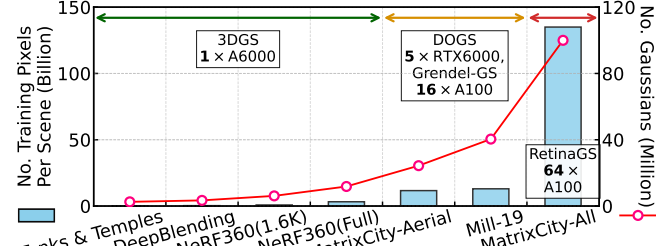


**Fig. 1: The trend of training 3D-GS for larger, more detailed scenes challenges both software and hardware. As the number of training pixels grows, software has evolved from 3DGS [13] to distributed frameworks like Grendel-GS [43] and RetinaGS [16], while hardware has scaled from 1 GPU to 64 GPUs, with Gaussians per scene increasing by 20-100×.**

(i.e., **retraining**) or post-process a pre-trained model through fine-tuning (i.e., **refinement**), both of which involve training and pose new challenges to applicability: Both retraining and fine-tuning require the full training pipeline, including access to often limited training data and computing resources, and can be time-consuming.

More recently, researchers aim to handle high-resolution (e.g., 1080p and 4K) and large-scale 3D reconstruction tasks with tens of millions of Gaussians [3, 16, 43]. However, given the huge number of Gaussians, 3D-GS cannot be trained nor rendered on a single GPU regardless of the batch size because it exceeds the single GPU memory capacity. As shown in Fig. 1, while small to moderate-scale scenes, like those in Tanks&Temples [14] and Mip-NeRF360 [2], can be processed on a single Nvidia A6000, large-scale scenes like MatrixCity [18] and Mill-19 [31] require distributed multi-GPU setups (e.g., 16-64) due to the additional memory demands of holding training images and optimizer states. The RUBBLE scene from Mill-19 [31], trained using Grendel-GS [43], contains over 40 million Gaussians and a 8.9 GB file size with a large training memory footprint. Unfortunately, existing 3D-GS compression work often overlook the compute requirements associated with high-resolution and large-scale 3D construction tasks, which in fact has a big impact to the deployment of 3D-GS with limited resources.

As an initial attempt to address these challenges, FCGS [4] introduces a generalizable model that compresses Gaussians without training, achieving state-of-the-art performance in under a minute. However, it falls short in two areas: (1) FCGS demands prohibitively high hardware resources, particularly memory (often over 24GiB), limiting its applicability and excluding resource-constrained mobile devices; (2) it lacks flexibility to adjust to rapidly varying compression needs, with each adjustment still taking about a minute.

To enable flexible and cost-efficient Gaussian compression, we introduce FlexGaussian, a training-free method with low compute requirements (e.g., mobile deployable) and instant reconfigurability

to adapt to varying compression needs. Specifically, unlike DNNs, where different layers exhibit diverse sensitivity, we observe that different Gaussian attribute channels exhibit different quantization sensitivity. As such, we propose to apply INT4/INT8 *Channel-wise Mixed-Precision Quantization (MPQ)* scheme with configurable bit-widths across Gaussian attribute channels. In addition, different from previous approaches for compressing 3D-GS that rely on learning K-means-based codebooks [25, 27], which require iterative optimizations, our approach is non-iterative, which significantly reduces the compression overhead while minimizing the compression error. We further show that such channel-wise sensitivity can be exploited for semi-structured pruning as well, and we introduce *Attribute-Discriminative Pruning (ADP)*, which enables fine-grained pruning at both the Gaussian and attribute granularities while eliminating the need to train the pruning mask [37], which greatly simplifies the pruning method for 3D-GS.
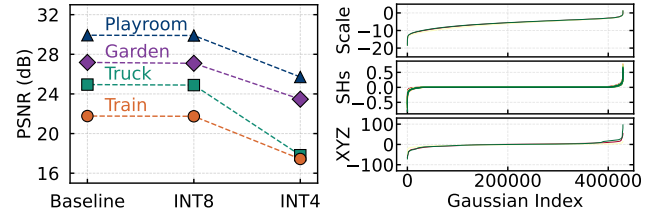
FlexGaussian exposes a broad optimization space to explore and identify the optimal compression rate within hardware-defined constraints. To navigate this space, we introduce *Fast Online Adaption (FOA)*, a lightweight yet effective adaptation module that jointly considers ADP and MPQ to identify the optimal parameter combinations for meeting hardware-specified constraints, such as target quality or compression ratio. With no training involved, the time cost to reach each trade-off spot is reduced to seconds, which in turn enables FOA to quickly search multiple spots.

Overall, in response to the rising costs of large-scale Gaussian training and varying compression needs, FlexGaussian eliminates the need for training in 3D-GS compression. This approach avoids the high-end hardware requirements incurred by storing training data and optimizer states on GPUs, reduces compression time costs, and offers flexibility in adjusting compression rates and qualities.

We conduct extensive experiments on both desktop and resource-constrained mobile platforms across diverse datasets to compress 3D-GS models for novel view synthesis, covering scene scales and image resolutions from small to large. Our evaluation shows that FlexGaussian achieves a significant compression ratio while maintaining high rendering quality, e.g., a 96.4% model size reduction with <1 dB drop in PSNR, with compression rates fully adjustable at negligible cost. Compared to state-of-the-art methods, FlexGaussian is 10-100× faster than training-involved methods and 1.7-2.1× faster than the training-free FCGS [4]. It also adapts rapidly to varying memory or bandwidth constraints, with each adjustment taking just 1 to 2 seconds, making it a flexible and cost-effective solution for Gaussian compression.

## 2 Related Work

**3D Gaussian compression**. Many studies focus on 3D-GS compression. EAGLES [9] introduces an encoder-decoder network to compress Gaussian attributes into a latency code. Scaffold-GS [22] uses anchor points to guide the distribution of 3D Gaussians and learns to predict their attributes at these points. LightGaussian [7] distills Spherical Harmonics(SHs) to a lower degree. Compressed3D [27] introduces sensitivity-aware clustering, quantization-aware fine-tuning, and entropy encoding to reduce 3D-GS size. Similarly, RDO-Gaussian [32] learns adaptive pruning of SHs and assigns varying



(a) RTN with varying bit-widths.  (b) Varying ranges in attributes.

**Fig. 2: Left: Quality of quantized 3D-GS with varying attribute bit-widths across scenes. Right: 3D-GS attributes show divergent ranges, with up to 200× difference between XYZ to SHs.**

SH degrees based on material and illumination needs. While promising, these methods all require significant training effort to identify redundancies and condense 3D Gaussian information. FCGS [4] develops a universal model for compression without training, but still suffers from insufficient memory even on GPUs like RTX 3090 and A6000. This hinders their applicability in resource-constrained scenarios, such as mobile devices or large-scale 3D-GS that requires expensive multi-GPU setups for training [16, 43]. Our work differs by introducing a flexible, cost-effective compression method that adaptively adjusts compression levels while maintaining comparable quality, with lower resource usage and suitability for even resource-constrained mobile devices.

**Slimmable scene representations.** Training a single model that allows for trade-offs between task quality and model size at runtime is crucial, especially given the varying resources available. This model can be classical Deep Neural Networks [12, 17, 38, 39] or recent scene representations such as NeRF [23] or 3D Gaussians [13] for rendering tasks. Recent research achieves slimmability in NeRF representations. CCNeRF [30] uses rank truncation to dynamically control model size and levels of detail through hybrid tensor rank decomposition. Similarly, SlimmeRF [40] also employs rank truncation but incrementally builds the rank during training. While there has been little prior work on achieving slimmability in 3D Gaussians, our work presents a hybrid method that dynamically reduces model size while maintaining quality, which cannot be accomplished by simply applying individual designs.

## 3 Challenges in Training-Free Compression

Post-training compression methods have demonstrated great efficiency in deep learning (DL) and large language model (LLM) compression, as they are applied without retraining or fine-tuning [8, 19, 36]. Among these strategies, two commonly used ones are quantization and pruning. However, there is few investigation into training-free compression for 3D-GS. Here, we briefly discuss the challenges of applying post-training quantization and pruning to 3D-GS.

**Non-negligible accuracy loss with ultra-low quantization.** For a 3D-GS model with $N \times M$ parameters, existing post-training quantization schemes [8, 19], such as Round-to-Nearest (RTN), reduce the storage by four times compared to FP32 by quantizing the parameter values to INT8. However, according to Fig. 2a, while 8-bit quantization largely maintains rendering quality across scenes, further reducing it to ultra-low bits, e.g., 4-bit, consistently lead to a sharp performance drop in rendered scenes.
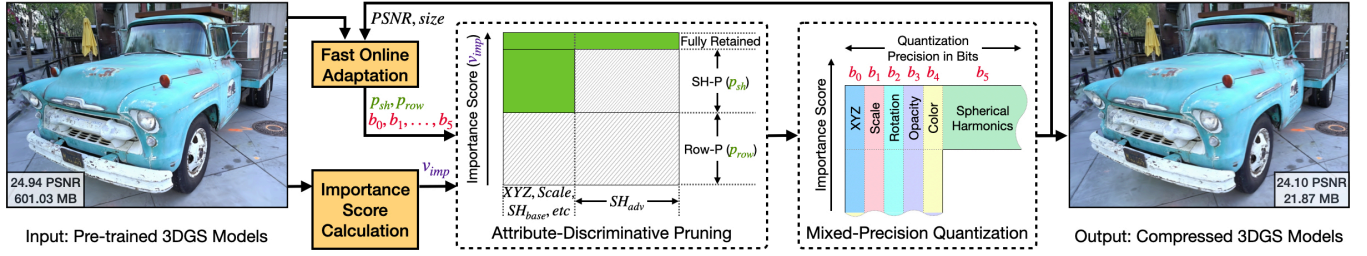
**Fig. 3: Overview of FlexGaussian. FlexGaussian first applies attribute-discriminative pruning (ADP) to obtain semi-structurally pruned Gaussian primitives. It then employs INT8/INT4 channel-wise mixed-precision quantization (MPQ) to further reduce the model size along the bit dimension. Finally, FlexGaussian introduces a novel lightweight online adaptation (FOA) algorithm to adaptively adjust the compression ratio for diverse scenes on hardware with different efficiency constraints.**

**Divergent attribute range.** To investigate why INT4 quantization leads to significant rendering quality drop, we analyze the per-attribute channel sensitivity both qualitatively and quantitatively. Fig. 2b shows that different attribute channels have dramatically different ranges (more details in § 4.2.1), and quantizing certain channels (e.g., scale) to ultra-low bits leads to large PSNR drop. *This large variance in the value range makes it difficult to use a fixed quantization range (usually the maximum value) for all channels while preserving rendering quality*, as the limited representation power for large-range channels leads to large quantization errors.

**Performance degradation from zero-shot Gaussian pruning.** Prior work prunes 3D-GS by removing the least significant 3D Gaussians based on a calculated importance score for each Gaussian. Common importance metrics include Opacity-Based Pruning [1, 42], which removes low-opacity Gaussian primitives due to their limited impact on visual quality; Magnitude-Based Pruning [20], which prunes primitives with fewer ray intersections as they contribute less to accuracy; and Volumetric-Based Pruning [29], which evaluates primitives by their 3D volume, with larger volumes indicating greater significance for rendering quality. And Global Significance Score [7], which combines Gaussian opacity, 3D volume, and pixel hit magnitude across training views to quantify its overall contribution to visual quality. However, prior work finds that directly pruning Gaussians (e.g., zero-shot pruning based on the magnitude of their opacity values or second-order approximation of the reconstruction error on the training views) leads to significant quality degradation [7, 10]. This is why existing methods typically employ retraining, fine-tuning, or knowledge distillation to adjust the remaining Gaussians and minimize quality loss [5, 7, 9, 15, 22, 25, 29, 32, 37]. However, as the compute requirement for scaling 3D-GS is growing, a lightweight and efficient method is desirable for compressing 3D-GS.

## 4 Methodology

### 4.1 Problem Formulation

The objective is to develop a lightweight and cost-effective method for training-free compression of 3D Gaussian models. Formally, the input of the problem includes an optimized 3D Gaussian model with a data layout of an $N \times M$ matrix, where each row ($N$) represents an individual 3D Gaussian, and each column ($M$) corresponds to an attribute channel. These channels include: ellipsoid's center coordinates (e.g., 3), scaling factor along each dimension (e.g., 3), rotation (encoded in quaternion: $w, x, y, z$), opacity, view-independent base

color in RGB ($SH_{base}$), and 3-degree spherical harmonics that enable view-dependent effects ($SH_{adv}$). Our goal is to obtain a *matching compressed 3D-GS model*, where its rendering quality is no less than $\epsilon$ (compression tolerance regime) in comparison with the uncompressed model. The solution should not require the training pipeline or access to the training images.

### 4.2 FlexGaussian Algorithm

FlexGaussian is flexible in compressing pre-trained 3D Gaussian models at negligible cost, to either achieve a target compression ratio or to meet a quality target. Fig. 3 highlights the key components of the system and demonstrates this procedure on the TRUCK scene from the Tanks&Temples dataset, achieving an 96.4% reduction in model size with <1 dB PSNR drop in just 20 seconds.

*4.2.1 INT4/INT8 Channel-wise Mixed-Precision Quantization (MPQ).*
Prior work on DNN compression find that different DNN layers exhibit diverse sensitivity [6, 33, 35], which motivates mixed-precision quantization, where more bits are assigned to sensitive layers to preserve model quality. Recent work has also explored spatial sensitivity in 3D-GS models by computing a second-order Hessian over a hypothetical perturbation to model weights and using the sensitivity score for pruning [10]. While being effective, it calculates the second order approximation for individual Gaussian, which is per-Gaussian and very expensive to calculate.

Intuitively, different attribute channels store distinct information (e.g., geometric, textual); do they exhibit varying quantization sensitivity? Our analysis in Fig. 4 shows the quality impact of applying 4-bit and 8-bit quantization to each attribute channel. The results show that channels in 3D-GS have varying sensitivities, and assigning the same number of bits to all channels is sub-optimal.

Based on the observation, we propose the *channel-wise mixed-precision quantization* method: We apply an ultra-low quantization bits (e.g., INT4) to each individual channel separately and keep the other channels in INT8 to find a normalized loss gap value between the uncompressed and quantized model. We use INT4 on channels with smaller gaps. While this approach does not explicitly consider the correlation between channels, we find that it is robust and efficient to select less sensitive columns. In addition, we observe minimal variation in channel sensitivity to bit-widths across different scenes, as certain Gaussian attributes (e.g., positions) are more prone to inaccuracies than others. Therefore, we use the same set of channel-specific bit-widths but can support the search procedure with minimal cost (several seconds) when necessary.
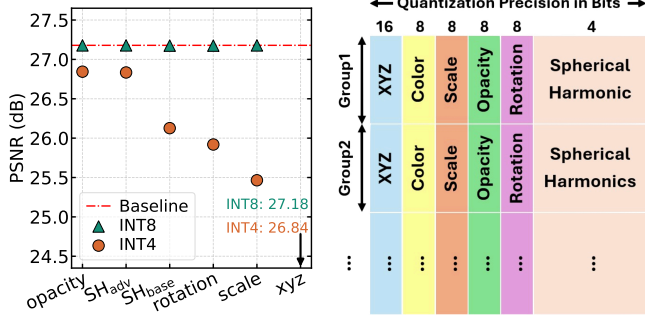
Fig. 4: Attribute channels exhibit divergent impact to rendering quality, evidencing the need for mixed-precision quantization that accounts for attribute sensitivity.
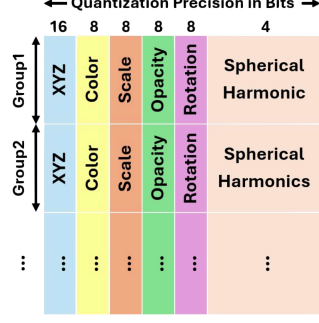


Fig. 5: Channel-wise mixed-precision quantization. The same color denotes the same bit-width. Group # indicates sub-groups with their own quantization ranges.
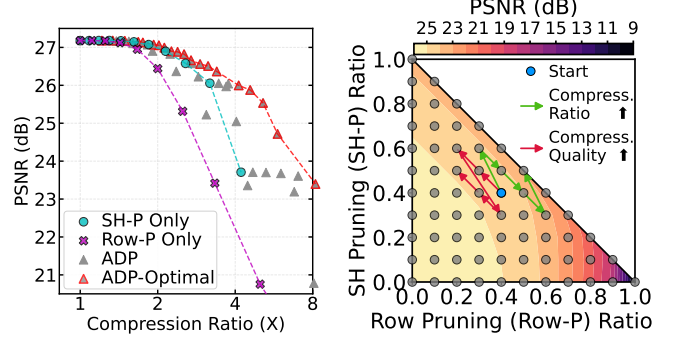


Fig. 6: A broad space for balancing quality and compression ratio, with optimal designs (ADP Optimal) outperforming pruning methods that ignore attribute sensitivity (Row-P and SH-P).



Fig. 7: FOA operates within the ADP design space, with optimal targets near the diagonal under quality constraints. Path with green arrows prioritizes compression, and red arrows favor quality.

While selective channel-wise mixed-precision quantization greatly reduces quantization error, we find that directly quantizing an entire channel with the same quantization range still leads to non-trivial accuracy degradation, as channels can contain millions of elements in some scenes. To tackle this, we use subchannel-wise grouped quantization for 3D-GS models. In each channel, we bucket sequential attributes together as sub-groups, e.g., all attributes of each set of channels are arranged to 1000 sub-groups. Each sub-group has its own quantization range, as shown in Fig. 5. Note that this method does not require iterative optimizations like prior methods [25, 27], which is hard to scale with increasing Gaussian primitives.

*4.2.2 Attribute-Discriminative Pruning (ADP).* Prior works effectively prune the number of Gaussian primitives but either overlook fine-grained pruning opportunities along the attribute dimension [7] or require iterative training to learn the pruning mask [37]. Our observation is that geometry attributes, such as position, rotation, and scale, have a greater impact on quality with a smaller data size (e.g., 11 vs. 59 in this case). In contrast, texture attributes such as $SH_{adv}$ require more storage (e.g., 45 attributes per Gaussian) but have less influence on quality (Fig. 4).

We exploit this opportunity by proposing Attribute-Discriminative Pruning, a semi-structured pruning method that effectively leverages the attribute-level importance properties to perform pruning. In this work, we consider a specific importance score from Light-Gaussian [7] – the global significant score of each Gaussian. In particular, the score $GS_n$ at the $n$−th row is calculated as:

$$GS_n = \sum_{n=1}^{MHW} \mathbb{1}(G(X_n), r_i) \cdot \sigma_n \cdot \gamma(\bigoplus_n) \qquad (1)$$

where $i$ represents a pixel, MHW represents the number of training views, image height, and width, respectively. $\mathbb{1}$ is the indicator function that determines whether a Gaussian intersects with a given ray $r_i$. $\sigma_n$ is the Gaussian's opacity, and $\gamma(\bigoplus_n)$ measure the 3D Gaussians' volume based on the scaling factors. The metric jointly considers how each Gaussian contributes to the global pixels, their opacity, and its volumetric density, which has been shown to be effective in pruning Gaussians. Note that the global significance score only needs to be calculated once given that the Gaussian attributes are fixed after training.

Instead of fully pruning or retaining Gaussian primitives, ADP enables partial pruning, where less critical attributes, such as $SH_{adv}$, can be partially pruned from non-critical Gaussians. Fig. 3 conceptually illustrates how ADP prunes Gaussians along both the primitive and attribute dimensions, with shaded areas representing the pruned parts. By sorting Gaussians based on importance scores in descending order, ADP retains the top $\alpha\%$ of Gaussians and discards the bottom $\beta\%$ that contribute minimally to quality. Unlike prior methods [7, 26, 37], ADP partially retains the remaining Gaussians $(1 - \alpha\% - \beta\%)$, preserving only critical attributes without re-training.

ADP has two parameters: the ratio of Gaussians pruned along the entire row (Row-P, i.e., $\alpha\%$) or partially for $SH_{adv}$ (SH-P, i.e., $1 - \alpha\% - \beta\%$). Their combination creates trade-offs between quality and compression ratio, as shown in Fig. 6. Notably, a Pareto-optimal frontier is formed with designs enabled by ADP, whereas simply applying Row-P as in prior work [7, 10] or SH-P individually leads to sub-optimal compression performance. ADP exposes a large set of parameter candidates, enabling flexible compression priorities.

*4.2.3 Fast Online Adaptation (FOA).* While identifying a single configuration for ADP and MPQ that generalizes across scenes is challenging, we find that the quality impact of MPQ is scene-insensitive, mainly due to its scene-agnostic operation, while the primary sensitivity arises from the importance score calculation for ADP. We extend the analysis in Fig. 6 to three additional scenes from different datasets and observe that, although the quality drop varies, the linear correlation between quality loss and pruning ratio remains consistent across scenes. Therefore, the Pareto-optimal frontier shares the same set of parameter pairs, narrowing down the optimal parameter configurations to a limited set of candidates.

We introduce a fast online adaptation (FOA) module that efficiently searches for the optimal design based on input scenes at negligible cost. Given a user-specified target quality or compression ratio, FOA searches in the candidate sets until reaches the best compression performance in observance to the constraints. While the search naturally can be done in parallel since each is independent to others, we exploit the fact that the Pareto-optimal frontier is convex, so that the priority trend at each direction is monotonious.
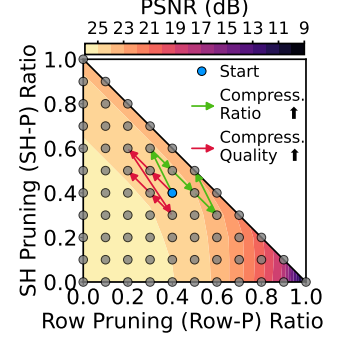
**Table 1: Quantitative comparison of quality scores, model sizes, and time costs for FlexGaussian and prior 3D-GS compression methods, measured on our desktop with a single RTX 3090 GPU. File size is in <u>MiB</u>, and time cost is in <u>Seconds</u>. Refinement-based, retraining-based, and training-free methods are shown in yellow , gray , and pink , respectively, with 3D-GS in blue as a reference. The best, second-best, and third-best methods in file size and time cost are shown in green , cyan , and orange .**

| Dataset Method \| Metric | Mip-NeRF360 | | | | | Tanks&Temples | | | | | Deep Blending | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | Size↓ | Time↓ | PSNR↑ | SSIM↑ | LPIPS↓ | Size↓ | Time↓ | PSNR↑ | SSIM↑ | LPIPS↓ | Size↓ | Time↓ |
| 3D-GS [13] | 27.21 | 0.815 | 0.214 | 795.26 | 1576.25 | 23.14 | 0.841 | 0.183 | 421.91 | 963.87 | 29.41 | 0.903 | 0.243 | 703.77 | 1654.11 |
| LightGaussian [7] | 26.96 | 0.800 | 0.244 | 52.05 | 917.61 | 23.14 | 0.818 | 0.222 | 27.87 | 426.18 | 28.82 | 0.891 | 0.273 | 45.85 | 686.04 |
| Compressed3D [27] | 27.02 | 0.803 | 0.236 | 28.81 | 272.84 | 23.34 | 0.836 | 0.192 | 17.25 | 191.45 | 29.44 | 0.902 | 0.250 | 25.28 | 247.98 |
| PUP 3DGS [10] | 26.66 | 0.789 | 0.267 | 79.53 | 427.99 | 22.46 | 0.799 | 0.248 | 42.19 | 189.69 | 29.40 | 0.903 | 0.256 | 70.38 | 270.76 |
| CompGS [26] | 26.99 | 0.801 | 0.250 | 21.08 | 2451.07 | 23.21 | 0.837 | 0.201 | 14.20 | 1397.70 | 29.98 | 0.911 | 0.250 | 15.15 | 1941.88 |
| FCGS-Raw [4]* | — | — | — | — | — | — | — | — | — | — | — | — | — | — | — |
| FCGS-Opt‡ | 27.04 | 0.799 | 0.231 | 62.34 | 53.57 | 23.36 | 0.839 | 0.186 | 31.12 | 30.09 | 29.61 | 0.902 | 0.243 | 55.41 | 42.61 |
| FlexGaussian (Ours) | 26.38 | 0.780 | 0.251 | 40.80 | 25.69 | 22.44 | 0.804 | 0.219 | 16.30 | 18.24 | 28.61 | 0.884 | 0.269 | 25.48 | 25.65 |

\* FCGS-Raw is the official implementation. It runs out of the 24 GiB VRAM of a single RTX 3090 on 7 of the 13 scenes. See per-scene results and discussions in Tbl. 3.

‡ FCGS-Opt is our customized version with reduced memory usage and faster speed, without compromising compression ratio or quality.

FOA begins with a one-time cost for model I/O and importance score calculation. At each search step, FOA duplicates the input Gaussians, then prunes, quantizes, dequantizes, and evaluates quality using standard metrics like PSNR, repeating until the best compression parameters are found. Each search yields different quality-size trade-offs within seconds, enabling instant reconfigurability, which is impractical for training-involved methods.

Fig. 7 illustrates the FOA procedure. The search aims to approach the diagonal (representing higher compression ratios) within the colorized zones (i.e., quality constraint). Starting from the initial state (blue), the search can traverse either along the green direction, prioritizing compression ratio, or the red direction, prioritizing quality. The traversal can occur bidirectionally when starting from the middle of the path, typically completing within a few seconds.

## 5 Evaluation

### 5.1 Evaluation Methodology

**Datasets.** We follow standard practices [7, 13] to evaluate Flex-Gaussian on three datasets: Mip-NeRF360 [2], Tanks&Temples [14], and Deep Blending [11], using 9, 2, and 2 scenes, respectively. These datasets span a wide variety of scenes, including both bounded indoor and unbounded outdoor scenarios, with diverse capture styles, object distributions, and levels of detail.

**Baselines.** We compare FlexGaussian with 3DGS [13] and several compression methods. We compare with FCGS [4], the only known training-free method recently introduced. For refinement-based baselines, we select LightGaussian [7], Compressed3D [27], and PUP-3DGS [10], and for retraining-based methods, we use CompGS [26], collectively referred to as **training-involved methods**. Our training-free design enables instant adaptation to various compression targets at minimal cost. In our main experiments, we target a <1 dB PSNR drop and evaluate adaptability separately.

Note that FCGS fails on 7 of the 13 test scenes due to excessive memory usage, exceeding the 24 GiB VRAM of the RTX 3090 used in our setup. We therefore compare with a customized variant, which reduces memory usage and time cost while preserving the original compression ratios and quality. Details are provided in Tbl. 3.

**Metrics.** We evaluate compression ratio, quality, and time cost. The compression ratio is calculated as the file size of the 3D-GS models divided by the uncompressed baseline. Compression quality is assessed using peak signal-to-noise ratio (PSNR), structural similarity (SSIM) [34], and perceptual similarity (LPIPS) [41], by comparing images rendered from models before and after compression. Time cost refers to the total time required for executing the compression methods, either the post-processing time for refinement-based methods or the total training time for retraining-based methods.

**Hardware setups.** Experiments are primarily conducted on a setup with an Intel Core i9-10900K, 64 GB of DRAM, and an Nvidia RTX 3090 GPU with 24 GB VRAM, unless otherwise noted. We also deploy our training-free design on the Nvidia Jetson Xavier [28], a mobile platform with 16 GB of shared memory between the CPU and a 512-core Volta GPU, to assess its lightweight advantage. To obtain large 3D-GS models, we train Grendel-GS [43] on 4 A100 GPUs, each with 40 GB of VRAM, interconnected via NVLink.

### 5.2 Main Results

**Quantitative results.** We evaluate the quality, time cost, and output file size of all methods in Tbl. 1. FlexGaussian achieves compression quality and file sizes comparable to training-involved methods, with a time cost of 1%-10%, typically under a minute. A similar trend holds against the FCGS variant, FCGS-Opt, which eliminates training but is still 1.7-2.1× slower than our method. Specifically, FlexGaussian reduces 94.9%, 96.1%, and 96.4% of the data across three datasets, with quality losses of 0.8 dB, 0.7 dB, and 0.8 dB, respectively — all well below the 1 dB PSNR drop constraint, and with an average time of under 30 seconds. Results for FlexGaussian on each scene are in Tbl. 3, along with a detailed comparison to FCGS and its variants. These varying compression ratios and quality losses show that FlexGaussian effectively identifies optimal compression parameters within the specified constraints. While training-involved methods offer a great balance of compression performance, they require pre-trained models and pose significant challenges, especially when training large-scale models.

Note that the time cost corresponds to the total execution time of the compression, including post-processing time for refinement or
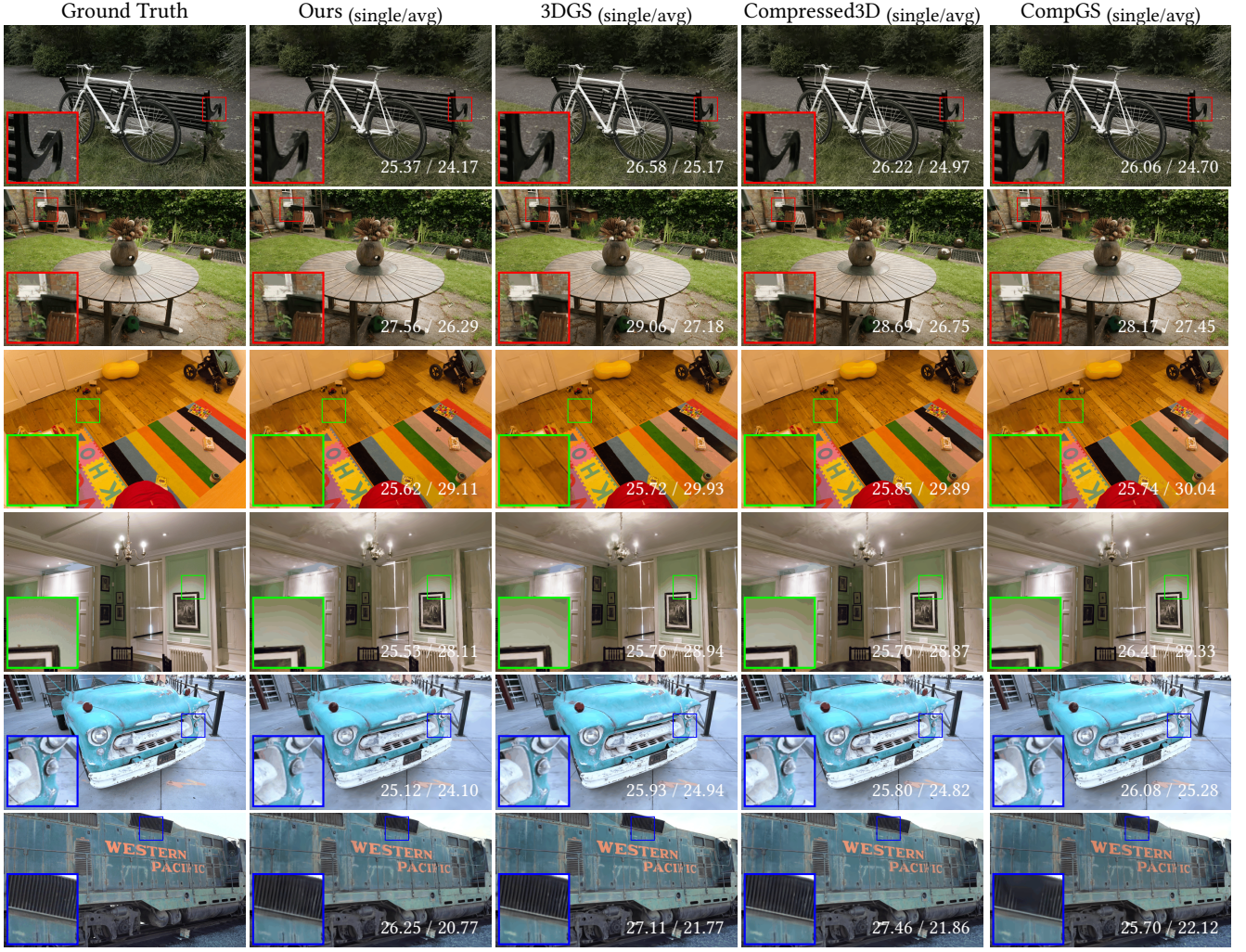
**Fig. 8: We compare the visual quality of novel views synthesized from models in FlexGaussian with the ground truth, 3D-GS [13], and two compression methods (refinement-based and retraining-based) using held-out test views. The evaluated scenes from top to bottom are BICYCLE from Mip-NeRF360, PLAYROOM from DeepBlending, and TRAIN from Tanks&Temples dataset. No obvious differences are observed between our method and the others. Results for other scenes are in the supplementary materials.**

training time for retraining approaches. For our method, the time includes data I/O, importance score computation, and FOA. This comparison might seem unfair to retraining-based methods, as all others start with pre-trained models. However, we note that Flex-Gaussian remains faster than all others, even when including the time required to train the 3D-GS models (3D-GS + FlexGaussian).

**Qualitative results.** FlexGaussian achieves visually indistinguishable results when comparing images rendered from uncompressed (3D-GS [13]) and compressed models (Fig. 8). We evaluate scenes from all three datasets, including indoor, outdoor, intricate, and large objects, with PSNR scores overlaid for the selected view and averaged across the scene. Additionally, we provide rendered videos in the supplementary materials that support these findings.

**Time breakdown.** We break down the total compression time of FlexGaussian, focusing on FOA 's online adaptation using the TRUCK scene in Fig. 9. Time measurement begins with the one-time cost of loading pre-trained 3D-GS models and camera parameters

to compute the importance score. Only complete evaluation sets are needed for quality loss calculation, while training camera poses suffice for importance computation. A baseline rendering is performed at the start of FOA for quality reference. Each FOA step deep copies the input model at the millisecond level before irreversible pruning, followed by quantization, dequantization, and quality evaluation, which together take about one second. The process repeats with updated parameters from a limited search set, all starting with the input model, until the optimal compression setup that satisfies the specified constraints is found. The full compression takes about 20.3 seconds, with 38.8%, 13.9%, 47.1%, and 0.2% of the time spent on data loading, importance score calculation, online adaptation, and model storage, respectively. Each adaptation step takes about 1.56 seconds on average, with pruning, quantization, dequantization, and rendering taking 0.22, 0.84, 0.32, and 0.18 seconds, respectively. The rendering time reflects rendering 32 views at approximately 178 FPS. Fast adaptation to varying compression trade-offs is crucial

for streaming Gaussians over fluctuating networks and on devices with diverse capabilities. Note that the times reported include the entire compression pipeline, with data I/O being a non-trivial part of our method but less significant for others.
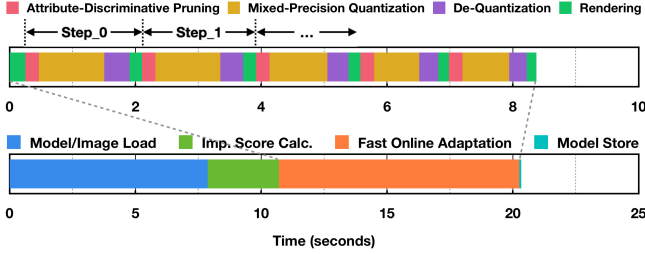


**Fig. 9: Fast online adaptation (FOA) adjusts quality and compression ratio every 2 seconds. The entire pipeline completes in under a minute, with data I/O being a non-trivial factor.**

**Hardware costs.** FlexGaussian 's elimination of training also reduces hardware requirements, making it arguably the first compression method **deployable on mobile platforms**, where both FCGS [4] (which runs out of 24 GiB VRAM even on a desktop setup, as shown by OOM in Tbl. 3) and training-involved methods are impractical. Fig. 10 shows the time breakdown for a mobile setup with the Nvidia Jetson Xavier, compared to a desktop setup with an RTX 3090. With only 512 CUDA cores in the VOLTA architecture and 12 GB of shared memory (4 GB of 16 occupied by the OS and other processes), our training-free design completes compression in minutes, detailed in Tbl. 3, still outperforming even the fastest Compressed3D [27] on the desktop. The increased time is mainly due to data I/O and importance score computation (∼7-24×), rather than FOA (∼3-5×), highlighting that Gaussian-intensive processing, which is heavily involved in training, is more demanding on less powerful hardware. Additionally, the extra memory required for training data and optimizer states, which far exceed the available memory on mobile devices, poses out-of-memory challenges when deploying training-involved methods on mobile platforms.

**Rate-distortion performance.** FlexGaussian enables fast exploration of the broad trade-off space between quality and compression ratio by adjusting pruning ratios for $Row-P$, $SH-P$, and bit-widths for attribute channels. While we set the quality constraint to <1 dB in our main experiments, the rate-distortion characteristics offer flexibility to optimize for either smaller file sizes or higher quality, spanning compression ratios from 8 to 256× and PSNR losses from 0.05 to 9 dB, as shown in Fig. 11 for the TRUCK scene. Each set of configurable parameters corresponds to a unique design point in the trade-off space and can be reconfigured to any other point in about a second (i.e., a step in Fig. 9). In contrast, most existing methods [7, 10, 26, 27] are limited to fixed design points and struggle to adapt to varying compression needs, while methods like RDO-Gaussian [32] and MesonGS [37] require retraining or refinement to switch between design points. Despite being training-free, FCGS [4] still requires a full feedforward pass (30-60 seconds) to reach a design point, while FOA in FlexGaussian handles such adjustments, with each taking only 1-2 seconds.
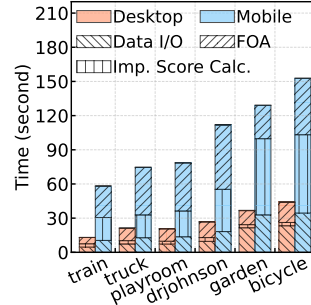


**Fig. 10: Time breakdown on desktop and mobile devices. Gaussian-intensive ops, like Imp. Score Calc., are more impacted by computing power.**
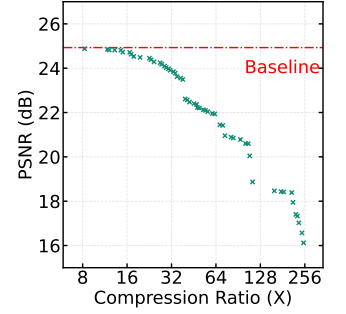


**Fig. 11: Rate-distortion performance of FlexGaussian, allowing instant access to compressions from 8 to 256× and PSNR losses of 0.05 to 9 dB.**

## 5.3 Analysis Results

**Large-scale 3D-GS.** FlexGaussian eliminates the computation and time costs associated with other training-involved methods, offering greater efficiency and applicability, particularly when training 3D-GS models requires specialized software or hardware. Large-scale Gaussian training with high-resolution data is increasingly used to represent larger scenes with high quality [21, 43], but it demands long training times, specialized software frameworks, and multi-GPU setups. We evaluate FlexGaussian 's applicability by compressing large models trained with Grendel-GS [43]. We train the BICYCLE scene at various resolutions—from its original (around 4K) to downsampled versions (2×, 4×, and 8×)—and at full resolution with checkpoints at different number of Gaussians, following Grendel-GS's densification approach. Tbl. 2 presents the time cost, peak memory usage, compression quality, and compression ratios for both training and applying FlexGaussian, all measured on a single Nvidia A100 GPU, with the quality drop restricted to <1 dB.

**Table 2: Training and compression costs vary with Gaussian scales, all measured under a quality constraint of <1 dB loss. High-resolution data and large Gaussians incur significant overhead in training and training-involved compression, evidencing the advantage of our training-free design.**

|  | Training Resolution | | | |
|---|---|---|---|---|
|  | 618×411 | 1237×822 | 2473×1643 | 4946×3286 |
| No. Gaussians (million) | 3.04 | 6.16 | 9.36 | 11.84 |
| Training Time (min.) | 17 | 33 | 73 | 190 |
| Training Mem. (GiB) | 6.81 | 14.00 | 23.54 | 39.37 |
| PSNR (dB) | 26.54 | 24.70 | 24.49 | 24.44 |
| FlexGaussian Time (sec.) | 20.08 | 38.64 | 64.71 | 104.31 |
| FlexGaussian Mem. (GiB) | 3.89 | 6.75 | 12.39 | 21.04 |
| PSNR (dB) | 25.64 | 24.10 | 23.72 | 23.89 |
| Compression Ratio (×) | 17.13 | 20.38 | 28.04 | 28.00 |

As image resolution and Gaussian counts increase, training costs in time and hardware rise sharply. This also challenges training-involved compression methods, which often require similar software and hardware resources. In contrast, FlexGaussian compresses Gaussian models at various scales with minimal resource and time consumption, all within the 1 dB quality constraint. This is made possible by eliminating the training process and the need for additional memory to store training data and optimizer states.

**Table 3: PSNR (dB), file size (MiB), and compression time (sec.) for each test scene.** $\boxed{\text{OOM}}$ : out of memory on RTX 3090 (24 GiB).

| Metrics | Method | Mip-NeRF360 | | | | | | | | | Tanks&Temples | | Deep Blending | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Bicycle | Bonsai | Counter | Flowers | Garden | Kitchen | Room | Stump | Treehill | Train | Truck | Drjohnson | Playroom |
| PSNR (dB) | 3D-GS [13] | 25.18 | 32.02 | 28.91 | 21.46 | 27.17 | 30.81 | 31.40 | 26.6 | 22.29 | 21.78 | 24.93 | 28.99 | 30.13 |
| | FCGS-Raw $1e^{-4}$ | OOM | 32.12 | 28.98 | OOM | OOM | 30.83 | 31.54 | OOM | OOM | 21.76 | OOM | OOM | 30.24 |
| | FCGS-Opt $1e^{-4}$ | 24.40 | 32.12 | 28.98 | 21.17 | 26.27 | 30.83 | 31.54 | 25.98 | 22.04 | 21.76 | 24.96 | 28.99 | 30.24 |
| | FCGS-Opt $16e^{-4}$ | OOM | 31.49 | 28.67 | OOM | OOM | 30.08 | 31.07 | OOM | 22.04 | 21.65 | OOM | OOM | 29.98 |
| | FlexGaussian | 24.18 | 31.05 | 28.24 | 20.57 | 26.30 | 30.03 | 30.67 | 25.72 | 21.29 | 20.83 | 24.16 | 28.13 | 29.29 |
| Size (MiB) | 3D-GS [13] | 1450.28 | 294.42 | 289.24 | 860.06 | 1379.99 | 438.10 | 376.85 | 1173.52 | 894.90 | 242.78 | 601.03 | 805.36 | 602.19 |
| | FCGS-Raw $1e^{-4}$ | OOM | 24.32 | 24.29 | OOM | OOM | 39.17 | 26.48 | OOM | OOM | 20.12 | OOM | OOM | 49.10 |
| | FCGS-Opt $1e^{-4}$ | 92.60 | 24.31 | 24.29 | 61.72 | 113.23 | 39.17 | 26.48 | 101.70 | 77.56 | 20.12 | 42.12 | 61.71 | 49.10 |
| | FCGS-Opt $16e^{-4}$ | OOM | 13.50 | 13.58 | OOM | OOM | 20.67 | 15.21 | OOM | 42.22 | 10.86 | OOM | 34.76 | 27.46 |
| | FlexGaussian | 70.97 | 24.56 | 19.82 | 31.20 | 89.80 | 36.70 | 22.03 | 57.99 | 23.98 | 10.73 | 21.91 | 29.19 | 22.02 |
| Time (Second) | FCGS-Raw $1e^{-4}$ Desktop | OOM | 20.75 | 20.81 | OOM | OOM | 29.73 | 24.41 | OOM | OOM | 19.83 | OOM | OOM | 37.34 |
| | FCGS-Opt $1e^{-4}$ Desktop | 92.60 | 20.30 | 20.29 | 61.72 | 86.83 | 28.39 | 23.85 | 84.14 | 64.04 | 18.69 | 41.48 | 48.22 | 36.99 |
| | FCGS-Opt $16e^{-4}$ Desktop | OOM | 13.15 | 12.94 | OOM | OOM | 18.65 | 16.36 | OOM | 40.60 | 11.20 | OOM | 33.72 | 25.18 |
| | FlexGaussian Desktop | 37.75 | 21.56 | 15.27 | 27.42 | 30.25 | 25.66 | 15.74 | 28.60 | 29.00 | 13.52 | 22.96 | 28.93 | 22.36 |
| | FlexGaussian Mobile | 135.49 | 94.37 | 68.34 | 91.33 | 117.18 | 122.48 | 82.22 | 66.31 | 82.21 | 59.42 | 73.88 | 113.45 | 77.15 |

**Comparison to FCGS.** FCGS [4] is a newly introduced Gaussian compression method that operates without refinement or retraining, making it highly relevant to our work. However, the off-the-shelf version fails on 7 out of 13 test scenes, requiring over 24 GiB of GPU VRAM—typical for GPUs like the RTX 3090 and A6000, as shown by the $\boxed{\text{OOM}}$ error in the FCGS-Raw in Tbl. 3, an inherent limitation of its design. To ensure a fair comparison across compression methods on the same hardware, we evaluate our enhanced variant, FCGS-Opt, with customized data partition strategies. FCGS-Opt reduces memory usage and execution time while maintaining the original compression quality and ratios. We follow FCGS's recommended setup, testing with two hyperparameters: $1e^{-4}$, prioritizing quality retention, and $16e^{-4}$, prioritizing smaller file sizes.

Tbl. 3 confirms that, under the same $1e^{-4}$ setting, FCGS-Opt reduces memory and time costs of FCGS-Raw at no compromise. However, it still fails due to out-of-memory errors when using the $16e^{-4}$ setting, which prioritize compression for smaller file sizes at the expense of quality. While FCGS offers flexibility in adjusting quality and file size trade-offs, its quality-focused configuration ($1e^{-4}$) lags behind FlexGaussian in both compression ratios and time costs, and its high compression ratio setting ($16e^{-4}$), which consumes excessive memory, limits its applicability—particularly on resource-constrained devices like mobile platforms. Additionally, each adjustment in FCGS requires a full feedforward pass that takes 30 to 60 seconds, whereas FlexGaussian completes the process in just 1 to 2 seconds with one step in FOA.

**Table 4: PSNR scores and model sizes for ablations on ADP and MPQ. Each contributes similarly and complements the other to maximize compression performance together.**

| PSNR (dB) | Size (MiB) | Garden | | Truck | | DrJohnson | |
|---|---|---|---|---|---|---|---|
| Baseline (3D-GS) [13] | | 27.18 | 1379.99 | 24.94 | 601.03 | 28.94 | 805.36 |
| + ADP | | 26.51 | 439.73 | 24.27 | 102.89 | 28.21 | 137.87 |
| + MPQ | | 26.81 | 216.65 | 24.69 | 93.97 | 28.81 | 124.43 |
| FlexGaussian (Ours) | | 26.29 | 89.58 | 24.10 | 21.87 | 28.11 | 28.98 |

**Ablation study.** We conduct ablation studies to assess the impact of each component. Our results show that ADP and MPQ must work **synergistically** to achieve the best compression performance, which is impractical when applying either individually in

a training-free setup. We analyze the contributions at the module level (i.e., ADP and MPQ) and then delve deeper into each module to examine the impact of various design decisions and strategies. Tbl. 4 compares the compressed models in rendering quality and file size, using one scene from each dataset across four configurations: the original 3D-GS [13] on our hardware setup, ADP individually, MPQ individually, and both combined. While MPQ achieves a slightly higher compression ratio and less quality loss than ADP, further reducing bit-widths in MPQ compromises quality significantly. Combining ADP with MPQ unlocks an additional 2.4 - 4.3× compression ratio while satisfying the specified quality constraints.

**Table 5: PSNR and model sizes for ablations on attribute-discriminative pruning, which boosts compression ratios while preserving quality over attribute-agnostic pruning.**

| PSNR (dB) | Size (MiB) | Garden | | Truck | | DrJohnson | |
|---|---|---|---|---|---|---|---|
| Baseline (3D-GS) [13] | | 27.18 | 1379.99 | 24.94 | 601.03 | 28.94 | 805.36 |
| Row-P | | 25.10 | 552.00 | 21.82 | 120.21 | 26.53 | 161.07 |
| SH-P | | 23.71 | 327.46 | 22.90 | 142.62 | 27.40 | 191.10 |
| ADP | | 26.51 | 439.73 | 24.27 | 102.89 | 28.21 | 137.87 |

**How is being aware of attribute sensitivity significant in pruning?** ADP prunes both Gaussian primitives and their attributes. Tbl. 5 compares pruning along each dimension (Row-P and SH-P) with ADP in terms of rendering quality (PSNR scores) and file sizes across scenes from three datasets, all with quantization deactivated. ADP clearly outperforms in preserving quality at similar compression ratios, as further supported by curves in Fig. 6.

**Table 6: Quality impact of subchannel-wise grouped quantization, halving quality loss at the same compression ratio.**

| | Garden | | | Bonsai | | |
|---|---|---|---|---|---|---|
| | $PSNR^{\uparrow}$ | $SSIM^{\uparrow}$ | $LPIPS^{\downarrow}$ | $PSNR^{\uparrow}$ | $SSIM^{\uparrow}$ | $LPIPS^{\downarrow}$ |
| Baseline (3D-GS) [13] | 27.18 | 0.861 | 0.115 | 31.98 | 0.938 | 0.208 |
| w/o Group Quant. | 26.45 | 0.832 | 0.148 | 30.25 | 0.915 | 0.237 |
| w/ Group Quant. | 26.81 | 0.846 | 0.131 | 31.19 | 0.928 | 0.219 |

**Does subchannel-wise grouped quantization in MPQ improve rendering quality? Yes.** Tbl. 6 compares the rendering quality of scenes from two datasets, both with and without grouped quantization, while keeping pruning deactivated. Subchannel-wise

grouped quantization significantly reduces quantization error, mitigating about half of the quality loss at the same compression ratio. We omit file sizes as this design choice has a negligible impact.

**How does importance score calculation impact compression performance?** While we use LightGaussian [7] for importance scores, FlexGaussian is independent of score calculation methods. To verify this, we test scores calculated from MesonGS [37] and compare rendering quality and compression ratios across scenes from three datasets in Tbl. 7. The results show that our system remains effective while meeting compression constraints. The performance gap indicates that more accurate importance estimation can improve both compression ratio and quality.

**Table 7: PSNR and model sizes of FlexGaussian with importance scores from MesonGS, demonstrating its compatibility.**

| PSNR (dB) \| Size (MiB) | Garden | | Truck | | DrJohnson | |
|---|---|---|---|---|---|---|
| Baseline (3D-GS) [13] | 27.18 | 1379.99 | 24.94 | 601.03 | 28.94 | 805.36 |
| MesonGS [37] | 26.30 | 80.40 | 23.99 | 21.87 | 28.24 | 28.98 |
| FlexGaussian (Ours) | 27.17 | 89.80 | 24.16 | 21.91 | 28.13 | 29.19 |

## 6 Conclusion and Future Work

We propose FlexGaussian, a flexible, cost-efficient, and training-free method for compressing 3D-GS models with high rendering quality, high compression ratios, and low adaption costs to varying compression demands. At its core are attribute-discriminative pruning, ultra-low bit mixed-precision quantization, and a fast online adaptation algorithm that supports dynamically adjusting the compression ratio. FlexGaussian achieves up to 96.4% compression with <1 dB PSNR drop in just seconds — two orders of magnitude faster than prior retraining- and refinement-based approaches. FlexGaussian is deployable on resource-constrained mobile devices, which is impractical for all other compression methods. For future work, we plan to explore further pushing the compression ratio by investigating ternary or even binary attributes as well as achieving even faster adaption speed via advanced online search algorithms.

## Acknowledgments

## References

[1] Muhammad Salman Ali, Maryam Qamar, Sung-Ho Bae, and Enzo Tartaglione. 2024. Trimming the fat: Efficient compression of 3d gaussian splats through pruning. *arXiv preprint arXiv:2406.18214* (2024).

[2] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. 2022. Mip-NeRF 360: Unbounded Anti-Aliased Neural Radiance Fields. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022.* IEEE, 5460–5469.

[3] Yu Chen and Gim Hee Lee. 2025. DOGS: Distributed-Oriented Gaussian Splatting for Large-Scale 3D Reconstruction Via Gaussian Consensus. *Advances in Neural Information Processing Systems* 37 (2025), 34487–34512.

[4] Yihang Chen, Qianyi Wu, Mengyao Li, Weiyao Lin, Mehrtash Harandi, and Jianfei Cai. 2024. Fast feedforward 3d gaussian splatting compression. *arXiv preprint arXiv:2410.08017* (2024).

[5] Yihang Chen, Qianyi Wu, Weiyao Lin, Mehrtash Harandi, and Jianfei Cai. 2025. Hac: Hash-grid assisted context for 3d gaussian splatting compression. In *European Conference on Computer Vision.* Springer, 422–438.

[6] Zhen Dong, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. 2019. Hawq: Hessian aware quantization of neural networks with mixed-precision. In *Proceedings of the IEEE/CVF international conference on computer vision.* 293–302.

[7] Zhiwen Fan, Kevin Wang, Kairun Wen, Zehao Zhu, Dejia Xu, and Zhangyang Wang. 2023. Lightgaussian: Unbounded 3d gaussian compression with 15x reduction and 200+ fps. *arXiv preprint arXiv:2311.17245* (2023).

[8] Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. 2022. GPTQ: Accurate Post-Training Quantization for Generative Pre-trained Transformers. *CoRR* abs/2210.17323 (2022).

[9] Sharath Girish, Kamal Gupta, and Abhinav Shrivastava. 2023. Eagles: Efficient accelerated 3d gaussians with lightweight encodings. *arXiv preprint arXiv:2312.04564* (2023).

[10] Alex Hanson, Allen Tu, Vasu Singla, Mayuka Jayawardhana, Matthias Zwicker, and Tom Goldstein. 2024. PUP 3D-GS: Principled Uncertainty Pruning for 3D Gaussian Splatting. *CoRR* abs/2406.10219 (2024).

[11] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. 2018. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (ToG)* 37, 6 (2018), 1–15.

[12] Liang Hou, Zehuan Yuan, Lei Huang, Huawei Shen, Xueqi Cheng, and Changhu Wang. 2021. Slimmable generative adversarial networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 7746–7753.

[13] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Trans. Graph.* 42, 4 (2023), 139–1.

[14] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. 2017. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)* 36, 4 (2017), 1–13.

[15] Joo Chan Lee, Daniel Rho, Xiangyu Sun, Jong Hwan Ko, and Eunbyung Park. 2024. Compact 3d gaussian representation for radiance field. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 21719–21728.

[16] Bingling Li, Shengyi Chen, Luchao Wang, Kaimin Liao, Sijie Yan, and Yuanjun Xiong. 2024. Retinags: Scalable training for dense scene rendering with billion-scale 3d gaussians. *arXiv preprint arXiv:2406.11836* (2024).

[17] Changlin Li, Guangrun Wang, Bing Wang, Xiaodan Liang, Zhihui Li, and Xiaojun Chang. 2021. Dynamic slimmable network. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition.* 8607–8617.

[18] Yixuan Li, Lihan Jiang, Linning Xu, Yuanbo Xiangli, Zhenzhi Wang, Dahua Lin, and Bo Dai. 2023. Matrixcity: A large-scale city dataset for city-scale neural rendering and beyond. In *Proceedings of the IEEE/CVF International Conference on Computer Vision.* 3205–3215.

[19] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. 2024. AWQ: Activation-aware Weight Quantization for On-Device LLM Compression and Acceleration. *Proceedings of Machine Learning and Systems* 6 (2024), 87–100.

[20] Weikai Lin, Yu Feng, and Yuhao Zhu. 2025. MetaSapiens: Real-Time Neural Rendering with Efficiency-Aware Pruning and Accelerated Foveated Rendering. In *Proceedings of the 30th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1.* 669–682.

[21] Tao Lu, Ankit Dhiman, R Srinath, Emre Arslan, Angela Xing, Yuanbo Xiangli, R Venkatesh Babu, and Srinath Sridhar. 2024. Turbo-GS: Accelerating 3D Gaussian Fitting for High-Quality Radiance Fields. *arXiv preprint arXiv:2412.13547* (2024).

[22] Tao Lu, Mulin Yu, Linning Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. 2024. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 20654–20664.

[23] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.

[24] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Trans. Graph.* 41, 4, Article 102 (July 2022), 15 pages. doi:10.1145/3528223.3530127

[25] KL Navaneet, Kossar Pourahmadi Meibodi, Soroush Abbasi Koohpayegani, and Hamed Pirsiavash. 2023. Compact3d: Compressing gaussian splat radiance field models with vector quantization. *arXiv preprint arXiv:2311.18159* (2023).

[26] KL Navaneet, Kossar Pourahmadi Meibodi, Soroush Abbasi Koohpayegani, and Hamed Pirsiavash. 2024. Compgs: Smaller and faster gaussian splatting with vector quantization. In *European Conference on Computer Vision.*

[27] Simon Niedermayr, Josef Stumpfegger, and Rüdiger Westermann. 2024. Compressed 3D Gaussian Splatting for Accelerated Novel View Synthesis. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024.* IEEE, 10349–10358.

[28] Nvidia. 2018. Jetson AGX Xavier Series, https://tinyurl.com/jetsonxavier. https://www.nvidia.com/en-sg/autonomous-machines/embedded-systems/jetson-agx-xavier/

[29] Panagiotis Papantonakis, Georgios Kopanas, Bernhard Kerbl, Alexandre Lanvin, and George Drettakis. 2024. Reducing the Memory Footprint of 3D Gaussian Splatting. *Proceedings of the ACM on Computer Graphics and Interactive Techniques* 7, 1 (2024), 1–17.

[30] Jiaxiang Tang, Xiaokang Chen, Jingbo Wang, and Gang Zeng. 2022. Compressible-composable nerf via rank-residual decomposition. *Advances in Neural Information Processing Systems* 35 (2022), 14798–14809.

[31] Haithem Turki, Deva Ramanan, and Mahadev Satyanarayanan. 2022. Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition.* 12922–12931.

[32] Henan Wang, Hanxin Zhu, Tianyu He, Runsen Feng, Jiajun Deng, Jiang Bian, and Zhibo Chen. 2024. End-to-End Rate-Distortion Optimized 3D Gaussian Representation. *arXiv preprint arXiv:2406.01597* (2024).

[33] Kuan Wang, Zhijian Liu, Yujun Lin, Ji Lin, and Song Han. 2019. Haq: Hardware-aware automated quantization with mixed precision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition.* 8612–8620.

[34] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* 13, 4 (2004), 600–612.

[35] Bichen Wu, Yanghan Wang, Peizhao Zhang, Yuandong Tian, Peter Vajda, and Kurt Keutzer. 2018. Mixed precision quantization of convnets via differentiable neural architecture search. *arXiv preprint arXiv:1812.00090* (2018).

[36] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. 2023. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning.* PMLR, 38087–38099.

[37] Shuzhao Xie, Weixiang Zhang, Chen Tang, Yunpeng Bai, Rongwei Lu, Shijia Ge, and Zhi Wang. 2025. MesonGS: Post-training Compression of 3D Gaussians via Efficient Attribute Transformation. In *European Conference on Computer Vision.* Springer, 434–452.

[38] Jiahui Yu and Thomas S Huang. 2019. Universally slimmable networks and improved training techniques. In *Proceedings of the IEEE/CVF international conference on computer vision.* 1803–1811.

[39] Jiahui Yu, Linjie Yang, Ning Xu, Jianchao Yang, and Thomas Huang. 2018. Slimmable neural networks. *arXiv preprint arXiv:1812.08928* (2018).

[40] Shiran Yuan and Hao Zhao. 2024. SlimmeRF: Slimmable Radiance Fields. In *2024 International Conference on 3D Vision (3DV).* IEEE, 64–74.

[41] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition.* 586–595.

[42] Zhaoliang Zhang, Tianchen Song, Yongjae Lee, Li Yang, Cheng Peng, Rama Chellappa, and Deliang Fan. 2024. LP-3DGS: Learning to Prune 3D Gaussian Splatting. *arXiv preprint arXiv:2405.18784* (2024).

[43] Hexu Zhao, Haoyang Weng, Daohan Lu, Ang Li, Jinyang Li, Aurojit Panda, and Saining Xie. 2024. On scaling up 3d gaussian splatting training. *arXiv preprint arXiv:2406.18533* (2024).