# 3DGS²: Near Second-order Converging 3D Gaussian Splatting

LEI LAN, University of Utah, USA

TIANJIA SHAO, Zhejiang University, China

ZIXUAN LU, University of Utah, USA

YU ZHANG, University of Utah, USA

CHENFANFU JIANG, University of California, Los Angeles, USA

YIN YANG, University of Utah, USA

**Fig. 1. Improve 3DGS training using parallelized local Newton.** We show 3DGS training can be substantially improved by switching from stochastic gradient descent to stochastic Newton. A Gaussian kernel has multiple attributes, and optimizing all the attributes of all the kernels as a whole leads to a very high-dimensional nonlinear problem. We leverage the fact that kernel attributes are weakly coupled and design a local Newton scheme to find the optimal value of each type of kernel parameter. We also exploit the spatial relation among input images to effectively mitigate the overshoot issue in stochastic optimization. As a result, our method is 5× to 10× faster than gradient-based 3DGS training. In the teaser, we show the reconstruction results using the-state-of-the-art 3DGS training algorithms, including vanilla 3DGS [Kerbl et al. 2023], AdR-Gaussian [Wang et al. 2024], EAGLES [Girish et al. 2025], Taming Gaussian [Mallick et al. 2024] as well as our method. Our method produces the result of the highest quality (in terms of SSIM, PSNR, and LPIPS) while using only one-tenth of iterations, making the reconstruction of complex 3D scenes in seconds instead of minutes.

3D Gaussian Splatting (3DGS) has emerged as a mainstream solution for novel view synthesis and 3D reconstruction. By explicitly encoding a 3D scene using a collection of Gaussian kernels, 3DGS achieves high-quality rendering with superior efficiency. As a learning-based approach, 3DGS training has been dealt with the standard stochastic gradient descent (SGD) method, which offers at most linear convergence. Consequently, training often requires tens of minutes, even with GPU acceleration. This paper introduces a (near) second-order convergent training algorithm for 3DGS, leveraging its unique properties. Our approach is inspired by two key observations. First, the attributes of a Gaussian kernel contribute independently to the image-space loss, which endorses isolated and local optimization algorithms. We exploit this by splitting the optimization at the level of individual kernel attributes, analytically constructing small-size Newton systems for each parameter group, and efficiently solving these systems on GPU threads. This achieves Newton-like convergence per training image without relying on the global Hessian. Second, kernels exhibit sparse and structured coupling across input images. This property allows us to effectively utilize spatial information to mitigate overshoot during stochastic training. Our method converges an order faster than standard GPU-based 3DGS training, requiring over 10× fewer iterations while maintaining or surpassing the quality of the compared with the SGD-based 3DGS reconstructions.

CCS Concepts: • **Computing methodologies → Rasterization**; **Reconstruction**; • **Mathematics of computing → Stochastic control and optimization**.

Additional Key Words and Phrases: 3DGS, Neural radiance fields, Stochastic Newton, Novel view synthesis, Optimization

## 1 INTRODUCTION

The reconstruction of 3D scenes is a fundamental challenge in computer graphics and vision. Recent advancements in learning-based techniques, such as Neural Radiance Fields (NeRFs) [Mildenhall et al. 2020], have re-banded this field. NeRF formulates the reconstruction process as a training task, encoding the color, texture, and geometry of 3D scenes through an implicit MLP network, achieving state-of-the-art results. A noteworthy follow-up is 3D Gaussian Splatting (3DGS) [Kerbl et al. 2023]. 3DGS represents a 3D scene explicitly with a collection of Gaussian splatting (GS) kernels. In contrast to NeRF, which relies on ray marching and volume rendering, 3DGS employs rasterization combined with a tile-splatting technique to generate novel views from unseen camera poses. This approach delivers superior runtime efficiency while still maintaining high-quality rendering results.

Just like other learning tasks, stochastic gradient descent (SGD) has been the mainstream modality for training 3DGS [Bottou 1998]. The GS kernels are initialized from SfM (structure-from-motion) [Snavely et al. 2006], and each kernel is assigned a position, a scaling vector, a rotation, an RGB color, and an opacity. The 3DGS training takes one input image as well as the corresponding camera pose.

Authors' addresses: Lei Lan, University of Utah, Salt Lake City, Utah, USA, lanlei.virhum@gmail.com; Tianjia Shao, Zhejiang University, Zhejiang, China, tjshao@zju.edu.cn; Zixuan Lu, University of Utah, Salt Lake City, Utah, USA, birdpeople1984@gmail.com; Yu Zhang, University of Utah, Salt Lake City, Utah, USA, zhan723284893@gmail.com; Chenfanfu Jiang, University of California, Los Angeles, Los Angeles, USA, Chenfanfu.Jiang@gmail.com; Yin Yang, University of Utah, Salt Lake City, Utah, USA, yangzzzy@gmail.com.

All the GS kernels are sorted based on the current camera depth, and an image-space loss value is evaluated between the rasterized image and the input image. The loss gradient is then used to update the configurations of all the GS kernels, and the learning rate is empirically adjusted at different stages of the training.

Such an SGD-based first-order optimization has been a default option for many deep-learning tasks, and it works well in general. Nevertheless, we argue that training 3DGS fundamentally differs from generic deep learning. Simply "copy-and-paste" the standard learning procedure for 3DGS is neither efficient nor effective. When a neural network is being trained, it is normally believed that parameters housed at neurons are *strongly two-way coupled*, meaning the perturbation at one parameter non-trivially influences all the other parameters as they are connected on the computational graph. This, however, is not the case for 3DGS training, where local greedy choice often suffices. For instance, in order to obtain a good scaling of a GS kernel to suppress the loss, one can always retrieve the best scaling along $x$ direction first, and adjust $y$- and $z$-direction scaling afterwards. In other words, they can be split into multiple sequential local greedy choices. Similarly, the color of a GS kernel does not contribute to a lot of pixels, suggesting color parameters are weakly coupled across GS kernels. Those observations suggest training the GS parameter with the global gradient may be sub-optimal.

In this paper, we propose a novel training algorithm for 3DGS. Our method is nearly second-order convergent based on local Newton's method. Concretely, we prioritize per-kernel parameters and calculate the second-order optimal update for each sub-group of parameters in series. The local updates are executed at all the pertaining GS kernels in parallel in a Jacobi-like manner. Because each local Newton only involves a small number of degrees of freedom (DOFs), we can analytically derive and assemble the local Hessian and gradient at a low cost for each GPU thread. The use of Newton's method avoids the necessity of line search, and the local optimizer always uses the default step size (e.g., the learning rate) of 1.0. Due to the weak coupling among DOFs, we observe strong second-order convergence during the training. We also deploy a novel sampling strategy, which effectively mitigates the overshoot issue of local optimizer for a given input image. We tested our method on multiple 3DGS datasets, and our algorithm is constantly 10× faster than SGD-based GPU 3DGS training with higher quality, making 3DGS training from minutes to seconds e.g. see Fig. 1.

## 2 RELATED WORK

Novel view synthesis from multi-view images/videos has been a core problem in computer graphics and vision, and a wide range of techniques have been developed for this purpose e.g., see [Aliev et al. 2020; Barron et al. 2021; Hedman et al. 2018; Kopanas et al. 2021; Lassner and Zollhöfer 2021; Lombardi et al. 2019, 2021; Thies et al. 2019; Zhang et al. 2022; Zhou et al. 2018]. NeRF [Mildenhall et al. 2020] implicitly encapsulates the information of the scene with an MLP net. By casting rays over this implicit representation, high-quality scene images can be synthesized. 3DGS [Kerbl et al. 2023], on the other hand, utilizes a set of GS kernels, which enable more efficient scene rendering via rasterization. They deliver state-of-the-art results and inspire many follow-up research efforts, such

as improving the rendering quality [Barron et al. 2022; Lu et al. 2024; Yu et al. 2024], increasing the scale of the scene [Kerbl et al. 2024; Song et al. 2025],

### 2.1 Accelerating NeRF

Although NeRF achieves high rendering quality, due to the expensive ray marching and the complexity of MLP, the training and rendering of NeRF are costly. Therefore, many methods have been proposed to accelerate NeRF's training and/or rendering. For example, Sun et al. [2022] proposed a dense voxel grid storing density and features, along with a post-activation interpolation on voxel density and robust optimization under several priors, leading to very fast convergence of NeRF optimization. Sparse voxel grids with learned features are also utilized to speed up the NeRF rendering [Hedman et al. 2021]. KiloNeRF [Reiser et al. 2021] utilizes thousands of small MLPs instead of a single large MLP to reduce the query cost, effectively accelerating the NeRF rendering. The octree data structure can be applied to speed up the rendering [Yu et al. 2021]. A GPU-friendly implementation [Chen et al. 2023] also helps NeRF achieve a faster rendering speed. Plenoxels [Fridovich-Keil et al. 2022] replaces MLPs with a sparse 3D grid with spherical harmonics. The state-of-the-art work is InstantNGP [Müller et al. 2022], which uses a hash grid and an occupancy grid to accelerate computation and a smaller MLP to represent density and appearance, and its training can be done in a few minutes. While these methods have had great success, due to the ray marching nature, the rendering performance of NeRF-based methods is, in general, slower than that of 3DGS.

### 2.2 Accelerating 3DGS

3DGS is the state-of-the-art radiance field representation that performs high-quality rendering in real time. Follow-up researches find that its rendering can be further accelerated due to high parameter counts and the redundant or uneven load in rasterization pipeline. For instance, EAGLES [Girish et al. 2025] utilized quantized embeddings to reduce per-point memory storage requirements and a coarse-to-fine training strategy to improve the training convergence speed, which largely reduces storage memory while speeding up the rendering. Hamdi et al. [2024] introduced Generalized Exponential Functions as a more memory-efficient alternative to Gaussians, which requires fewer particles to represent the scene, reducing the memory storage and increasing the rendering speed. LightGaussian [Fan et al. 2024] prunes Gaussians with minimal impact on visual quality, condenses spherical harmonic coefficients, and applies vector quantization, which substantially compresses Gaussians and largely boosts the rendering speed. C3DGS [Lee et al. 2024] proposes a learnable mask to remove Gaussians that have minimal impact on overall quality, also reducing the storage and increasing the rendering speed. Speedy-Splat [Hanson et al. 2024] precisely localizes Gaussians by computing a tight bounding box around their extent and uses soft and hard pruning to effectively reduce the Gaussian numbers and largely speed up the rendering time as well as training time. AdR-Gaussian [Wang et al. 2024] early culls Gaussian-Tile pairs with low splatting opacity, and proposes a balancing algorithm for pixel thread load, achieving a large increase of rendering and training speed. Other examples include changing the

densification heuristics to reduce the number of Gaussians [Fang and Wang 2024; Kheradmand et al. 2024; Rota Bulò et al. 2025] and improving the underlying differentiable rasterizer [Feng et al. 2024; Ye et al. 2024]. Their goals are mainly compressing Gaussians and accelerating the Gaussian rendering, with the increased training speed as a byproduct. Taming 3DGS [Mallick et al. 2024] utilizes a new densification strategy to reduce Gaussians, together with the backpropagation with per-splat parallelization, to speed up the Gaussian training on limited resources. 3DGS-LM [Höllein et al. 2024] accelerates the training convergence by replacing the stochastic gradient descent (SGD) optimizer like ADAM [Kingma 2014] by Levenberg-Marquardt optimizer. The training process converges 30% faster. Nevertheless, the Gauss-Newton optimizer is a pseudo-second-order method. It does not utilize the weak coupling feature of the training. As a result, our method outperforms 3DGS-LM by a significant margin.

## 3 BACKGROUND

3DGS represents a radiance field using a set of Gaussian kernels $\mathcal{K}$. A GS kernel $k$ is parameterized with a set of trainable attributes, including the center position $\boldsymbol{p}_k \in \mathbb{R}^3$, the opacity $\sigma_k \in \mathbb{R}$, and the covariance matrix $\boldsymbol{A}_k \in \mathbb{R}^{3 \times 3}$. $\boldsymbol{A}_k$ is a positive semi-definite definite matrix which only has six independent DOFs. A common practice is to decompose $\boldsymbol{A}_k$ as $\boldsymbol{A}_k = \boldsymbol{R}_k \boldsymbol{S}_k \boldsymbol{S}_k^\top \boldsymbol{R}_k^\top$ with a diagonal scaling matrix $\boldsymbol{S}_k = \text{diag}(s_{k,x}, s_{k,y}, s_{k,z})$ and a rotation matrix $\boldsymbol{R}_k$. A more compact representation is to use a unit quaternion $\boldsymbol{q}_k \in \mathbb{R}^4$, $\|\boldsymbol{q}_k\| = 1$ as well as a scaling vector $\boldsymbol{s}_k = [s_{k,x}, s_{k,y}, s_{k,z}]^\top$ to encode the orientation and the geometry of the $k$-th Gaussian ellipsoid. A kernel also has its own color information, which is expressed with spherical harmonics (SH) coefficient vector $\boldsymbol{c}_k$. To render an image of the scene, GS kernels are sorted based on the current camera depth and projected onto the image plane. The color of an image pixel at $(m, n)$ is computed as:

$$\boldsymbol{c}(m, n) = \sum_{k \in \mathcal{K}} G_k(m, n) \sigma_k \tilde{\boldsymbol{c}}_k(\boldsymbol{r}_k, \boldsymbol{c}_k) \prod_{j=1}^{k-1} \left(1 - G_j(m, n) \sigma_j\right). \quad (1)$$

Here, $G_k(m, n)$ denotes the 2D Gaussian weight of the $k$-th kernel at the pixel whose image coordinate is $(m, n)$. $\boldsymbol{r}_k$ stands for the view direction from the camera to the kernel center. $\tilde{\boldsymbol{c}}(\boldsymbol{r}_k, \boldsymbol{c}_k)$ is the view-dependent color information expressed via spherical harmonics (SH). A loss function $L$ measures the difference between the rasterized image and the ground truth image, which consists of a per-pixel color difference and an SSIM loss.

The training procedure of the vanilla 3DGS searches for a better setup of all the GS kernels by $\boldsymbol{x} \leftarrow \alpha \Delta \boldsymbol{x}$. Here, $\boldsymbol{x}$ is a vector of $|\mathcal{K}| \cdot 14$ dimensions, concatenating parameters of all the $|\mathcal{K}|$ kernels. At each iteration, an incremental improvement $\Delta \boldsymbol{x}$ is chosen as the negative gradient of the loss $\Delta \boldsymbol{x} = -\frac{\partial L}{\partial \boldsymbol{x}}$, and the step size $0 < \alpha < 1$ is used to prevent overshooting. The ground-truth loss is defined over all the input images. Because calculating the global gradient is costly, the training uses a sampled gradient based on each input image and progressively reduces the loss in an image-by-image manner.

The above procedure is well-known as a standard way for non-linear optimization of a high-dimension and nonlinear problem i.e.,
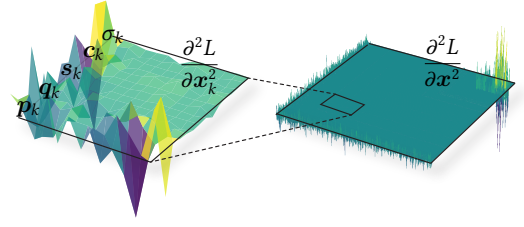


Fig. 2. **Hessian visualization.** We plot the values of local and global Hessian matrices for one training image. The local Hessian (left) is for all the parameters of a kernel $\boldsymbol{x}_k = [\boldsymbol{p}_k^\top, \boldsymbol{q}_k^\top, \boldsymbol{s}_k^\top, \boldsymbol{c}_k^\top, \sigma_k]^\top$, while the global Hessian (right) is for all the kernels' parameters. The variation of the matrices across different DOFs suggests (very) weak coupling amount parameters.

$\min_{\boldsymbol{x}} L(\boldsymbol{x})$ using stochastic gradient descent (SGD). SGD is commonly considered linearly convergent as it leverages the first-order Taylor expansion of the target loss [Nocedal and Wright 1999]. Therefore, it is not surprising to see that 3DGS training needs tens of thousands of iterations to adequately reduce the training loss.

The convergence of the training procedure can be substantially improved using Newton's method, which quadratically expands $L(\boldsymbol{x}^*)$ at $\boldsymbol{x}$ such that:

$$L(\boldsymbol{x}^*) = L^* = L + \boldsymbol{g}^\top \Delta \boldsymbol{x} + \frac{1}{2} \Delta \boldsymbol{x}^\top \boldsymbol{H} \Delta \boldsymbol{x} + O(\|\Delta \boldsymbol{x}\|^3). \quad (2)$$

Here, $\boldsymbol{x}^*$ is a *local minimizer* of $L$ around current $\boldsymbol{x}$. $\boldsymbol{g} = \left(\frac{\partial L}{\partial \boldsymbol{x}}\right)^\top$, and $\boldsymbol{H} = \frac{\partial^2 L}{\partial \boldsymbol{x}^2}$ are the gradient and Hessian of the loss $L$. Ignoring the cubic error term, the corresponding DOF update $\Delta \boldsymbol{x}$ can be computed from:

$$\Delta \boldsymbol{x} = \arg \min_{\boldsymbol{y}} L + \boldsymbol{g}^\top \boldsymbol{y} + \frac{1}{2} \boldsymbol{y}^\top \boldsymbol{H} \boldsymbol{y}$$

via solving the linear system of:

$$\boldsymbol{H}(\boldsymbol{x}) \Delta \boldsymbol{x} = -\boldsymbol{g}(\boldsymbol{x}). \quad (3)$$

Newton's method possesses many desired properties for nonlinear optimization. First, if $O(\|\Delta \boldsymbol{x}\|^3)$ is reasonably small, Newton's method converges quadratically. It normally needs much fewer iterations than GD. More importantly, Newton's method does not need the line search [Nocedal and Wright 1999], i.e., $\alpha$ by default is one. When the quadratic approximation of $L^*$ is valid, and $\Delta \boldsymbol{x}$ is second-order optimal.

## 4 OUR METHOD

Unfortunately, solving the global Newton system at each iteration is computationally expensive, if not impossible. Even the assembly of the matrix is costly and requires tedious implementation. As a result, Newton's method never appears as a viable candidate for 3DGS training. Instead of resorting to the global Newton procedure, we argue that local Newton optimization is also efficient and effective.

This reasoning is drawn based on several important observations. First of all, 3DGS training never aims for a global optimum. Just as other deep learning tasks, obsessing over whether a global optimum is achieved holds little practical importance. Second, training DOFs at a GS kernel are weakly coupled. Numerically, this can

be confirmed by visualizing the full Hessian of 3DGS training. As shown in Fig. 2, there exist sparsely localized "spikes" of the values of $\frac{\partial^2 L}{\partial x^2}$. This implies the influence among kernel parameters is not global. Solving the full Newton system does not offer a proportional improvement in convergence relative to its computational cost. Similarly, the Hessian of the attributes of one GS kernel is also concentrated, and positional DOFs play a more important role compared with parameters related to color. Departing from generic learning tasks, the training data for 3DGS is highly structured. It is possible for us to exploit the spatial relation among input images to effectively mitigate the oscillation (e.g., overshoot) during the sampled optimization. In the following subsections, we elaborate on the details of the proposed training algorithm.

## 4.1 Training data preparation

Given a set of input training images $\mathcal{T} = \{I^1, I^2, \cdots, I^{|\mathcal{T}|}\}$, we kick off the training following the vanilla 3DGS by computing a set of initial GS kernels via SfM points. We then estimate a scene center and an encapsulating bounding sphere, which contains all the GS kernels. We project camera poses to the surface of this bounding sphere. The distance between a pair of input images is then defined as the spherical distance between their corresponding camera projections. With the distance defined, we can then obtain the $K$ nearest neighbors (KNNs) of a given input image $I^t$, which intuitively are images that share similar camera perspectives of $I^t$. We call $I^t$ i.e., the current training image, the *primary target*, and the set of its KNNs *secondary targets*, denoted as $\mathcal{N}_t$. Training images that are distant from $I^t$ i.e., $\mathcal{T} - (\{I^t\} + \mathcal{N}_t)$, are considered less relevant as only a small fraction of local kernels are visible from them. On the other hand, secondary targets are more tightly coupled with $I^t$. In our implementation, we set $K = 3$ and noticed this value offers a good balance between convergence and efficiency. As to be discussed later, Hessians and gradients of secondary targets are sparsely evaluated, which accumulate to the primary Hessian/gradient to prevent overshoot (i.e., see Sec. 4.3).
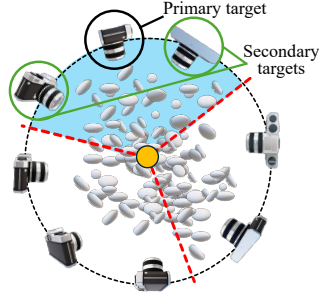


Fig. 3. **Data preparation.** We construct a bounding sphere of the scene containing all the GS kernels. Camera poses corresponding to input training images are projected to the sphere surface, and the spherical distance between them is used as the metric to identify KNNs for training images. The current training image is the primary target while its KNNs are the secondary targets.

## 4.2 Stochastic Local Newton

As mentioned, each GS kernel $k$ houses the position $\boldsymbol{p}_k$, opacity $\sigma_k$, orientation $\boldsymbol{q}_k$, scaling $\boldsymbol{s}_k$, as well as a SH color vector $\boldsymbol{c}_k$. Rather than learning the full parameter of $\boldsymbol{x}_k = [\boldsymbol{p}_k^\top, \sigma_k, \boldsymbol{q}_k^\top, \boldsymbol{s}_k^\top, \boldsymbol{c}_k^\top]^\top$ as one single variable, we split the optimization into multiple sub-steps, following the order of position, orientation, scaling, opacity, and

color. We note that position is of the most importance and should always be optimized first. Shape parameters (orientation, scaling) are prioritized over the color information. This intuition is straightforward to follow: without the right placement of a kernel, the correlated pixels are blank, leaving no opportunities for color-wise adjustment. However, we do not have a strong preference regarding the order of orientation and scaling or color and opacity. Numerically speaking, they become nearly decoupled after the position has been determined.

Each parameter of kernel $k$ is optimized via a local Newton solve given a certain training input $I^t$:

$$\Delta \boldsymbol{y}_k = -\left(\frac{\partial^2 L^t}{\partial \boldsymbol{y}_k^2}\right)^{-1} \left(\frac{\partial L^t}{\partial \boldsymbol{y}_k}\right)^\top, \tag{4}$$

where $\boldsymbol{y}_k$ refers to a specific set of variables to be optimized. The loss function depends on the RGB values at a pixel $(m, n)$ from 3DGS rasterization $\boldsymbol{c}(m, n)$ and from the training input $I^t$ i.e., $\boldsymbol{c}^t(m, n)$. This allows us to use the chain rule to compute the Hessian and gradient of the loss:

$$\frac{\partial L^t}{\partial \boldsymbol{y}_k} = \frac{\partial L^t}{\partial \boldsymbol{c}} \frac{\partial \boldsymbol{c}}{\partial \boldsymbol{y}_k}, \quad \frac{\partial^2 L^t}{\partial \boldsymbol{y}_k^2} = \left(\frac{\partial \boldsymbol{c}}{\partial \boldsymbol{y}_k}\right)^\top \frac{\partial^2 L^t}{\partial \boldsymbol{c}^2} \frac{\partial \boldsymbol{c}}{\partial \boldsymbol{y}_k} + \frac{\partial L^t}{\partial \boldsymbol{c}} \cdot \frac{\partial^2 \boldsymbol{c}}{\partial \boldsymbol{y}_k^2}. \tag{5}$$

In our implementation, the loss function is the superposition of the L2 loss and SSIM loss such that $L^t = L_2^t + \lambda L_S^t$. To assemble Eq. (4) for all the kernel attributes, we first compute the first and second derivatives of L2 loss and SSIM loss w.r.t. the rasterized color:

$$\frac{\partial L_2^t}{\partial \boldsymbol{c}(m, n)} = \frac{1}{3|I^t|} \left(\boldsymbol{c}(m, n) - \boldsymbol{c}^t(m, n)\right), \quad \frac{\partial^2 L_2^t}{\partial \boldsymbol{c}^2(m, n)} = \frac{\boldsymbol{I}_3}{3|I^t|}, \tag{6}$$

$$\frac{\partial L_S^t}{\partial \boldsymbol{c}(m, n)} = \frac{1}{3|I^t|} \sum_{j=0}^{H} \sum_{i=0}^{W} \frac{\partial f_S\left(m + i - \frac{W}{2}, n + j - \frac{H}{2}\right)}{\partial \boldsymbol{c}(m, n)}, \tag{7}$$

$$\frac{\partial^2 L_S^t}{\partial \boldsymbol{c}^2(m, n)} = \frac{1}{3|I^t|} \sum_{j=0}^{H} \sum_{i=0}^{W} \frac{\partial^2 f_S\left(m + i - \frac{W}{2}, n + j - \frac{H}{2}\right)}{\partial \boldsymbol{c}^2(m, n)}. \tag{8}$$

Here, $|I^t|$ denotes the total number of pixels in the input training $I^t$. $\boldsymbol{I}_3$ is a $3 \times 3$ identity matrix. $F_S(a, b)$ measures the SSIM error between $\boldsymbol{c}(a, b)$ and $\boldsymbol{c}^t(a, b)$ at pixel $(a, b)$ [Wang et al. 2004] as the sum of errors of RGB components i.e., $F_S(a, b) = F_{S,R} + F_{S,G} + F_{S,B}$. The error of each color component e.g., the red component $F_{S,R}$ is:

$$F_{S,R}(a, b) = \frac{(2\mu\mu^t + C_1)(2\tilde{\rho}^t + C_2)}{\left((\mu\mu^t)^2 + C_1\right)\left((\rho\rho^t)^2 + C_2\right)}, \tag{9}$$

where $\mu$, $\mu^t$ and $\rho$, $\rho^t$ are respectively the weighted means and variances of the red values within a $W \times H$ Gaussian filter centered at pixel $(a, b)$ in the rasterized image and training image $I^t$. $\tilde{\rho}^t$ is the weighted covariance. $C_1 = 0.01^2$ and $C_2 = 0.03^2$ are two constants. SSIM uses a filter of $W = H = 11$. $F_{S,G}$ and $F_{S,B}$ are calculated similarly.

It is clear that $F_{S,R}$ only depends on the red color value, and it is irrelevant to the green or blue components. Therefore, we have:

$$\frac{\partial F_S(a,b)}{\partial \boldsymbol{c}(m,n)} = \frac{\partial F_S}{\partial [c_R, c_G, c_B]^\top} = \left[\frac{\partial F_{S,R}}{c_R}, \frac{\partial F_{S,G}}{c_G}, \frac{\partial F_{S,B}}{c_B}\right], \quad (10)$$

$$\frac{\partial^2 F_S(a,b)}{\partial \boldsymbol{c}^2(m,n)} = \text{diag}\left(\frac{\partial^2 F_{S,R}}{c_R^2}, \frac{\partial^2 F_{S,G}}{c_G^2}, \frac{\partial^2 F_{S,B}}{c_B^2}\right). \quad (11)$$

The first and second derivatives of the red SSIM error w.r.t. the red component (and other color components) can be computed as:

$$\frac{\partial F_{S,R}}{\partial c_R} = \frac{f_1}{f_2 f_3} g_0 + \frac{f_0}{f_2 f_3} g_1 - \frac{f_0 f_1}{f_2^2 f_3} g_2 - \frac{f_0 f_1}{f_2 f_3^2} g_3, \quad (12)$$

$$\frac{\partial^2 F_{S,R}}{\partial^2 c_R} = \left(\frac{g_1}{f_2 f_3} - \frac{f_2 g_3 + f_3 g_2}{f_2^2 f_3^2} f_1\right) g_0^2 + \left(\frac{g_0}{f_2 f_3} - \frac{f_2 g_3 + f_3 g_2}{f_2^2 f_3^2} f_0\right) g_1^2$$

$$+ \left(\frac{-(f_0 g_1 + f_1 g_0)}{f_2^2 f_3} + \frac{2 f_2 f_3 g_2 + f_2^2 g_3}{(f_2^2 f_3)^2} f_0 f_1\right) g_2^2$$

$$+ \left(\frac{-(f_0 g_1 + f_1 g_0)}{f_2 f_3^2} + \frac{2 f_2 f_3 g_3 + f_3^2 g_2}{(f_2 f_3^2)^2} f_0 f_1\right) g_3^2, \quad (13)$$

where

$$f_0 = 2\mu\mu^t + C_1, \quad f_1 = 2\tilde{\rho}^t + C_2, \quad f_3 = (\rho\rho^t)^2 + C_2,$$

$$g_0 = 2 w_{i,j} \mu^t, \quad g_1 = 2 w_{i,j}(\boldsymbol{c}(m,n) - \mu^t),$$

$$g_2 = 2 w_{i,j} \mu, \quad g_3 = 2 w_{i,j}(\boldsymbol{c}(m,n) - \mu).$$

Here, $w_{i,j} = w(m - a + \frac{H}{2}, n - b + \frac{W}{2})$ is the filter's weight.

The per-attribute Hessian and gradient can then be computed with Eq. (5), where the first and second derivatives of $\boldsymbol{c}(m,n)$ i.e., how the pixel colors change w.r.t. kernel attributes are evaluated at each individual attribute.

*Position solve.* The position of the kernel $k$ is not full-rank. During the optimization, we only permit planner position adjustment of $k$ that is perpendicular to the camera's look at $r_k$. To extra the actual DOFs of the kernel position, we build a basis matrix $\boldsymbol{U}_k = [\boldsymbol{u}_x, \boldsymbol{u}_y] \in \mathbb{R}^{3 \times 2}$ such that:

$$\boldsymbol{u}_y = \frac{1}{\boldsymbol{r}_k \cdot \boldsymbol{r}_k}(\boldsymbol{r}_k \otimes \boldsymbol{r}_k)[0,1,0]^\top, \quad \boldsymbol{u}_x = \frac{1}{\sqrt{\boldsymbol{r}_k \cdot \boldsymbol{r}_k}} \boldsymbol{r}_k \times \boldsymbol{u}_y. \quad (14)$$

The update of the kernel position is then re-parameterized with $\boldsymbol{v}_k \in \mathbb{R}^2$:

$$\Delta \boldsymbol{p}_k = \boldsymbol{U}_k \Delta \boldsymbol{v}_k. \quad (15)$$

We first compute the partial derivatives w.r.t. the full positional DOFs via:

$$\frac{\partial \boldsymbol{c}}{\partial \boldsymbol{p}_k} = \frac{\partial \boldsymbol{c}}{\partial \tilde{\boldsymbol{c}}_k}\frac{\partial \tilde{\boldsymbol{c}}_k}{\partial \boldsymbol{p}_k} + \frac{\partial \boldsymbol{c}}{\partial G_k}\left(\frac{\partial G_k}{\partial \pi_k}\frac{\partial \pi_k}{\partial \boldsymbol{p}_k} + \frac{\partial G_k}{\partial \Sigma_k}\frac{\partial \Sigma_k}{\partial \boldsymbol{p}_k}\right), \quad (16)$$

$$\frac{\partial^2 \boldsymbol{c}}{\partial \boldsymbol{p}_k^2} = \frac{\partial \boldsymbol{c}}{\partial \tilde{\boldsymbol{c}}_k} \cdot \frac{\partial^2 \tilde{\boldsymbol{c}}_k}{\partial \boldsymbol{p}_k^2} + \frac{\partial \boldsymbol{c}}{\partial G_k}\left(\frac{\partial G_k}{\partial \pi_k} \cdot \frac{\partial^2 \pi_k}{\partial \boldsymbol{p}_k^2} + \frac{\partial G_k}{\partial \Sigma_k} : \frac{\partial^2 \Sigma_k}{\partial \boldsymbol{p}_k^2}\right.$$

$$+ \left(\frac{\partial \pi_k}{\partial \boldsymbol{p}_k}\right)^\top \frac{\partial^2 G_k}{\partial \pi_k^2}\frac{\partial \pi_k}{\partial \boldsymbol{p}_k} + \left(\frac{\partial \Sigma_k}{\partial \boldsymbol{p}_k}\right)^\top : \frac{\partial^2 G_k}{\partial \Sigma_k^2} : \frac{\partial \Sigma_k}{\partial \boldsymbol{p}_k}$$

$$\left.+ 2\left(\frac{\partial \Sigma_k}{\partial \boldsymbol{p}_k}\right)^\top : \frac{\partial^2 G_k}{\partial \pi_k \partial \Sigma_k} : \frac{\partial \pi_k}{\partial \boldsymbol{p}_k}\right). \quad (17)$$

The Hessian and the gradient of the generalized kernel position $\boldsymbol{v}_k$ can then be computed as:

$$\frac{\partial \boldsymbol{c}}{\partial \boldsymbol{v}_k} = \boldsymbol{U}_k^\top \frac{\partial \boldsymbol{c}}{\partial \boldsymbol{v}_k}, \quad \frac{\partial^2 \boldsymbol{c}}{\partial \boldsymbol{v}_k^2} = \boldsymbol{U}_k^\top \frac{\partial^2 \boldsymbol{c}}{\partial \boldsymbol{p}_k^2} \boldsymbol{U}_k. \quad (18)$$

Here, $\pi_k = \pi_k(\boldsymbol{p}_k) \in \mathbb{R}^2$ is the image space projection of the kernel center $\boldsymbol{p}_k$ a.k.a. the normalized device coordinate. $\Sigma_k \in \mathbb{R}^{2 \times 2}$ is 2D covariance matrix. $\tilde{\boldsymbol{c}}_k = \tilde{\boldsymbol{c}}_k(\boldsymbol{r}_k)$ is view-dependent SH color. The first and second derivatives of those intermediate variables are given in the supplementary document.

*Scaling solve.* A potential risk in optimizing scaling and orientation is redundancy. It is the projected 2D Gaussian ellipse that directly influences the pixel color on the rasterized image. Different combinations of scaling and orientation can result in the same projection. This brings the nullspace to the Hessian and undermines the convergence. To avoid this issue, we follow the strategy used for position solve, and perform the local optimization in a reduced space.

Specifically, to optimize $\boldsymbol{s}_k$, we only concern how it would affect the eigenvalue of the 2D covariance matrix $\Sigma_k$ i.e., the lengths of two axes of the projected ellipse. By eigenvalue decomposition $\Sigma_k = \boldsymbol{V}_k^\top \Lambda_k \boldsymbol{V}_k$, we have:

$$\Lambda_k = (\boldsymbol{V}_k \boldsymbol{J}_k \boldsymbol{W}_k \boldsymbol{R}_k) \text{diag}\left(s_{k,x}^2, s_{k,y}^2, s_{k,z}^2\right)(\boldsymbol{J}_k \boldsymbol{W}_k \boldsymbol{R}_k \boldsymbol{V}_k)^\top, \quad (19)$$

where $\Lambda_k \in \mathbb{R}^{2 \times 2} = \text{diag}(\lambda_{min}, \lambda_{max})$ is the diagonal matrix of eigenvalues. $\boldsymbol{V}_k \in \mathbb{R}^{2 \times 3}$ includes the corresponding eigenvectors. $\boldsymbol{W}_k$ and $\boldsymbol{J}_k$ represent the view matrix and the Jacobian matrix of orthographic projection $\boldsymbol{P}$. Since the updates of $s_{k,x}$, $s_{k,y}$, and $s_{k,z}$ only alter $\Lambda_k$, $\boldsymbol{V}_k$ remains unchanged. $\boldsymbol{J}_k$, $\boldsymbol{W}_k$ and $\boldsymbol{R}_k$ are also constant. Hence, we have:

$$(\boldsymbol{V}_k \boldsymbol{J}_k \boldsymbol{W}_k \boldsymbol{R}_k) \text{diag}\left(s_{k,x}, s_{k,y}, s_{k,z}\right) = \text{diag}(\lambda_{min}, \lambda_{max}), \quad (20)$$

leading to

$$\underbrace{\boldsymbol{E}_2^\top : (\boldsymbol{V}_k \boldsymbol{J}_k \boldsymbol{W}_k \boldsymbol{R}_k) : \boldsymbol{E}_3}_{T_k} \boldsymbol{s}_k = \boldsymbol{\lambda}_k = \left[\begin{array}{c} \lambda_{min} \\ \lambda_{max} \end{array}\right], \quad (21)$$

where $\boldsymbol{E}_2$ and $\boldsymbol{E}_3$ are 3-order tensors padding a vector to a diagonal matrix such that $\Lambda_k = \boldsymbol{E}_2 \cdot \boldsymbol{\lambda}_k$ and $\boldsymbol{S}_k = \boldsymbol{E}_3 \cdot \boldsymbol{s}_k$.

The gradient and Hessian of the full scaling factor $\boldsymbol{s}_k$ are:

$$\frac{\partial \boldsymbol{c}}{\partial \boldsymbol{s}_k} = \frac{\partial \boldsymbol{c}}{\partial G_k}\frac{\partial G_k}{\partial \Sigma_k} : \frac{\partial \Sigma_k}{\partial \boldsymbol{\lambda}_k}\frac{\partial \boldsymbol{\lambda}_k}{\partial \boldsymbol{s}_k},$$

$$\frac{\partial^2 \boldsymbol{c}}{\partial \boldsymbol{s}_k^2} = \frac{\partial \boldsymbol{c}}{\partial G_k}\left(\left(\frac{\partial \Sigma_k}{\partial \boldsymbol{\lambda}_k}\frac{\partial \boldsymbol{\lambda}_k}{\partial \boldsymbol{s}_k}\right)^\top \frac{\partial^2 G_k}{\partial \Sigma_k^2} : \frac{\partial \Sigma_k}{\partial \boldsymbol{\lambda}_k}\frac{\partial \boldsymbol{\lambda}_k}{\partial \boldsymbol{s}_k} + \frac{\partial G_k}{\partial \Sigma_k} : \frac{\partial \Sigma_k}{\partial \boldsymbol{\lambda}_k}\frac{\partial^2 \boldsymbol{\lambda}_k}{\partial \boldsymbol{s}_k^2}\right). \quad (22)$$

They are then projected into the least-square subspace of $T_k$ as:

$$\frac{\partial \boldsymbol{c}}{\partial \boldsymbol{\lambda}_k} = T_k(T_k^\top T_k)^{-\top}\frac{\partial \boldsymbol{c}}{\partial \boldsymbol{s}_k}, \quad \frac{\partial^2 \boldsymbol{c}}{\partial \boldsymbol{\lambda}_k^2} = T_k(T_k^\top T_k)^{-\top}\frac{\partial^2 \boldsymbol{c}}{\partial \boldsymbol{s}_k^2}(T_k^\top T_k)^{-1}T_k^\top. \quad (23)$$

*Rotation solve.* We optimize the rotation $q_k$ by constraining its axis to align with $r_k$ so the rotation is orthogonal to the scaling. This constraint allows us to re-write the rotation update as $\Delta q_k = [\cos \theta_k, \cos \theta_k r_k]^\top$, where $\theta_k$ is the rotation angle to be optimized:

$$\frac{\partial c}{\partial \theta_k} = \frac{\partial c}{\partial G_k} \frac{\partial G_k}{\partial \Sigma_k} \frac{\partial \Sigma_k}{\partial \theta_k},$$

$$\frac{\partial^2 c}{\partial \theta_k^2} = \frac{\partial c}{\partial G_k} \left( \frac{\partial \Sigma_k}{\partial \theta_k}^\top \frac{\partial^2 G_k}{\partial \Sigma_k^2} : \frac{\partial \Sigma_k}{\partial \theta_k} + \frac{\partial G_k}{\partial \Sigma_k} : \frac{\partial^2 \Sigma_k}{\partial \theta_k^2} \right). \quad (24)$$

Instead of the linear update $q \leftarrow q + \Delta q$, the rotation update is done with:

$$q \leftarrow \Delta q \cdot q. \quad (25)$$

*Opacity solve.* The opacity $\sigma_k$ is a scalar variable. However, its value is not arbitrary as $\sigma_k \in [0, 1)$, suggesting inequality constraints need to be imposed to avoid out-out-boundary values. We use the interior-point method [Potra and Wright 2000] by adding two logarithmic barrier terms into the local loss such that:

$$L^t \leftarrow L^t - \alpha_\sigma \ln \sigma_k + \ln(1 - \sigma_k). \quad (26)$$

The gradient of the opacity is:

$$\frac{\partial c}{\partial \sigma_k} = G_k \left( \prod_{j=1}^{k-1} \left(1 - G_j \sigma_j\right) \right) \left( c_k - \sum_{i=k+1} G_i \sigma_i c_i \prod_{j=k+1}^{i-1} \left(1 - G_j \sigma_j\right) \right), \quad (27)$$

and it has a vanished second derivative e.g., $\frac{\partial^2 c}{\partial \sigma_k^2} = 0$. The Hessian and gradient w.r.t. $L^t$ should also incorporate barrier loss i.e., $-\frac{1}{\sigma_k} - \frac{1}{1-\sigma_k}$ and $\frac{1}{(1-\sigma_k)^2} - \frac{1}{\sigma_k^2}$.

*Color solve.* The color of the kernel is encoded with SH bases $B_k$, which also depends on the current position of the kernel, i.e., $B_k(r_k(p_k))$. Because $p_k$ is always the first to be optimized, we decouple this relation during the color solve so that $r_k$ is now constant, and $\tilde{c}_k = B_k c_k$ is a linear function of the SH vector $c_k$.

Similar to $F_S$, the optimization of $c_k$ can be simplified by decoupling each color component e.g., between $c_R$ and $c_{k,R}$ as:

$$\frac{\partial c_R}{\partial c_{k,R}} = \sum_{k \in \mathcal{K}} G_k \sigma_k \prod_{j=1}^{k-1} \left(1 - G_j \sigma_j\right) B_{k,R}, \quad (28)$$

where $B_{k,R}$ is a row vector corresponding to the red component in $B_k$. It is easy to see that the second derivative $\frac{\partial^2 c_R}{\partial c_{k,R}^2} = 0$ is vanished.

### 4.3 Overshoot

A common challenge for stochastic optimization is *overshoot*. Due to the large amount of training images, it is impractical to evaluate the ground-truth Hessian $H^\star$ and gradient $g^\star$ because the loss is accumulated over the entire training set $L = L^1 + L^2 + \cdots + L^{|\mathcal{T}|}$. When we optimize the attributes of GS kernels under the view of the current training image $I^t$, the optimization is not aware of how the result could influence the loss values in other images. If the per-image optimization is over-aggressive, the reduction of $L^t$ can get outweighed by increases at other images, and the global loss $L$ becomes actually worsened. We refer to this numerical phenomenon as overshoot, which severely undermines the global convergence.

Preventing overshoot is expensive in general, as the global information of the loss function is needed. As a result, damped batch-based optimization is often employed [Kingma 2014; Qian 1999] together with the hyper-parameter tuning of the learning rate. Because of the second-order convergence, overshoot is more likely to occur without the constraint of the learning rate. However, if a learning rate is still necessary, why should we use Newton's method at all?

We give an efficient solution to this dilemma of stochastic training of 3DGS. Given a training image $I^t$ and the current GS kernel $k$ being trained, we define its *complementary* loss $\bar{L}^t$ such that:

$$\bar{L}^t = \sum_{j \in \mathcal{T}, j \neq t} L^j, \text{ for } \frac{\partial L^j}{\partial x_k} \neq 0. \quad (29)$$

In other words, $\bar{L}^t$ includes all the losses that are relevant to the kernel's attribute. Ideally, if the local training is in the form of $\min_{x_k}(L^t + \bar{L}^t)$, overshoot can be fully avoided. Unfortunately, doing so needs to traverse all the training data, which is clearly prohibitive. Given the spatial arrangement of training images, we note that $\bar{L}^t$ is not evenly distributed over $\mathcal{T}$. Rather, $\bar{L}^t$ concentrates on images sharing similar camera perspectives of $I^t$, i.e., the secondary targets (as mentioned in Sec. 4.1). Therefore, we approximate $\bar{L}^t$ as the total image losses of the secondary targets $\mathcal{N}_t$ and re-form the per-kernel training to the following format:

$$\min_{y_k} \left( L^t + \tilde{L}^t \right), \text{ where } \tilde{L}^t = \sum_{j \in \mathcal{N}_t} L^j. \quad (30)$$

The addition of $\tilde{L}^t$ is intended to prevent overshoot so that $\Delta y_k$ does not accidentally increase the loss in $\mathcal{N}_t$. We do not expect $\Delta y_k$ to play a primary role in reducing $\tilde{L}^t$ as this is handled by the respective local Newton solve. The actual gradient and Hessian of $\frac{\partial \tilde{L}^t}{\partial y_k}$ and $\frac{\partial^2 \tilde{L}^t}{\partial y_k^2}$ are of less interest as they act more like a self-adjusting pre-conditioner of $\Delta y_k$. Therefore, we downsample images in $\mathcal{N}_t$ for improved training efficiency. The gradient and Hessian of $\tilde{L}^t$ are computed in a similar way as discussed in the previous subsection.

The experiment is consistent with our analysis (see Sec. 5.2). The loss from the secondary targets offers an effective way to mitigate overshoot at a marginal computational cost. Solving Eq. (30), the global loss converges stably and hardly shows any oscillations at different training batches without the learning-rate-based step size trimming.

## 5 EXPERIMENTAL RESULTS

We implemented our pipeline on a desktop computer with an intel i7-12700 CPU and an Nvidia 3090 RTX GPU. Our algorithm was implemented using C++ and CUDA. We used standard datasets and metrics as in 3DGS [Kerbl et al. 2023]. Specifically, the datasets comprise all scenes from Mip-NeRF360 [Barron et al. 2022], two scenes from Tanks & Temples [Knapitsch et al. 2017], and two scenes from DeepBlending [Hedman et al. 2018]. Our experiments primarily measure the training convergence, time performance, and the quality of the resulting reconstruction. In general, our training method is 5× to 10× faster than the vanilla GS training using one order fewer iterations. The executable is also available in the supplemental material.

A representative example is reported in Fig. 4. In the figure, we plot the convergence curves of 3DGS training of six different scenes given the input image. Being gradient-based, vanilla 3DGS training converges much slower than our method. It is often the case that one iteration using the proposed per-attribute Newton lowers the loss more effectively than 100 gradient descent iterations. Because local Hessian and gradient are small-size, the computational time for each iteration is only 20% – 25% slower compared with one gradient descent.

### 5.1 Optimization order

An important strategy allowing local optimization is the decoupling of the kernel attribute. To further validate this, we plot convergence curves of three representative 3DGS reconstruction cases under different orders for per-attribute training. The result is reported in Fig. 5. It can be seen that optimizing the positions and geometries of GS kernels should be carried out before color information. Once kernel positions have been determined, orders of scaling and orientation do not matter since they become decoupled. The order between opacity and color is also of less importance. If we choose to optimize the color information first, on the other hand, training also converges but often to a lower-quality local minimum.

### 5.2 Trade off of secondary target

Another major contributor to our excellent convergence is mitigating overshoot by estimate $\tilde{L}$ at secondary target $\mathcal{N}_t$. A trade-off we have to decide is the size of $\mathcal{N}_t$ and the convergence rate. Incorporating more training into $\mathcal{N}_t$ improves the convergence but at the cost of more expensive local optimization since the Hessian and gradient of training data in the secondary target are also needed. To this end, we select three classic training scenarios and compare the training time and iteration counts with different $|\mathcal{N}_t|$ i.e., $|\mathcal{N}_t|= 0$, $|\mathcal{N}_t|= 3$, $|\mathcal{N}_t|= 8$. It can be clearly seen from Fig. 6 that increasing $|\mathcal{N}_t|$ from 0 to a small quantity effectively improves the convergence. Further enlarging $\mathcal{N}_t$ to 8 is able to better resolve this issue at the cost of more expensive computation. This is consistent with our previous analysis and endorses the design of our training algorithm.

### 5.3 Comparison with existing methods

While many efforts have been investigated to improve the rendering of a reconstructed 3DGS scene, there are relatively fewer works focusing on 3DGS training. We compared our method with vanilla 3DGS [Kerbl et al. 2023], AdR-Gaussian [Wang et al. 2024], EAGLES [Girish et al. 2025], 3DGS-LM [Höllein et al. 2024], and Taming 3DGS [Mallick et al. 2024]. The detailed benchmark is reported in Tab. 1. AdR-Gaussian [Wang et al. 2024] and EAGLES [Girish et al. 2025] focus on improving the efficiency of the rasterization of GS kernels and the memory footprint. Therefore, the total 3DGS training time may also be shortened. Our method is orthogonal to AdR-Gaussian or EAGLES, as we focus on the underlying numerical method for 3DGS training. Combining our method with faster rasterization techniques leads to even faster training. Nevertheless, we report our training performance using the standard 3DGS rasterization pipeline to avoid confusion. 3DGS-LM [Höllein et al. 2024] replaces SGD with global LM (Levenberg-Marquardt)

optimization [Ranganathan 2004] aiming for improved training convergence. Unfortunately, global LM optimization is not able to effectively extract localized nonlinearity (i.e., as shown in Fig. 2). Our method outperforms 3DGS-LM by a significant margin. Taming 3DGS [Mallick et al. 2024] aims to improve the training efficiency by using fewer GS kernels. It is more effective on hardware platforms with limited computing resources. With our experiment setup, Taming 3DGS performs similarly as AdR-Gaussian. It remains a first-order training modality and is slower than our method.

## 6 CONCLUSION & LIMITATION

This paper introduces a second-order convergent training method for 3DGS. We are inspired by the numerical properties of the loss Hessian and switch SGD-based linear optimization to local-Newton-based raining. The decoupled DOFs allow us to analytically derive the local Hessian for each type of kernel attribute and solve the resulting system efficiently in parallel. We also exploit the spatial correlation among input images to approximate the down-sampled global loss so that per-batch 3DGS training does not overshoot. This method outperforms existing method by a significant margin. It demonstrates a strong second-order convergent behavior and only needs one-tenth iterations. As a result, it reduced the total training time by one order without accuracy compromise.

We prove the feasibility of the second-order training for 3DGS reconstruction. While our results are encouraging, there are still many limitations that could be improved in the future. As the vanilla 3DGS, our algorithm is only tested for static scenes. While reconstructing dynamic scenes may be possible, moving GS kernels generates dynamically varying coupling, which may negatively impact the training effectiveness. Evaluating the Hessian of the kernel position is expensive and stands as one major bottleneck of our current pipeline. It may be possible to replace local Newton with Gauss-Newton (e.g. as in [Höllein et al. 2024] but at local DOFs) or quasi-Newton methods for different kernel attributes to further improve the efficiency. It is also of interest for us to explore auxiliary spatial data structures such as AABB or spatial hashing to speed up the kernel sorting [Lefebvre and Hoppe 2006].

## REFERENCES

Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. 2020. Neural Point-Based Graphics. In *Computer Vision − ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII* (Glasgow, United Kingdom). Springer-Verlag, Berlin, Heidelberg, 696–712.

Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. 2021. Mip-NeRF: A Multiscale Representation for Anti-Aliasing Neural Radiance Fields. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. 5835–5844.

Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. 2022. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5470–5479.

Léon Bottou. 1998. Online algorithms and stochastic approximations. *Online learning in neural networks* (1998).

Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. 2023. MobileNeRF: Exploiting the Polygon Rasterization Pipeline for Efficient Neural Field Rendering on Mobile Architectures. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 16569–16578.

Zhiwen Fan, Kevin Wang, Kairun Wen, Zehao Zhu, Dejia Xu, and Zhangyang Wang. 2024. LightGaussian: Unbounded 3D Gaussian Compression with 15x Reduction and 200+ FPS. arXiv:2311.17245 [cs.CV]

Guangchi Fang and Bing Wang. 2024. Mini-Splatting: Representing Scenes with a Constrained Number of Gaussians. In *Computer Vision − ECCV 2024*, Aleš Leonardis,
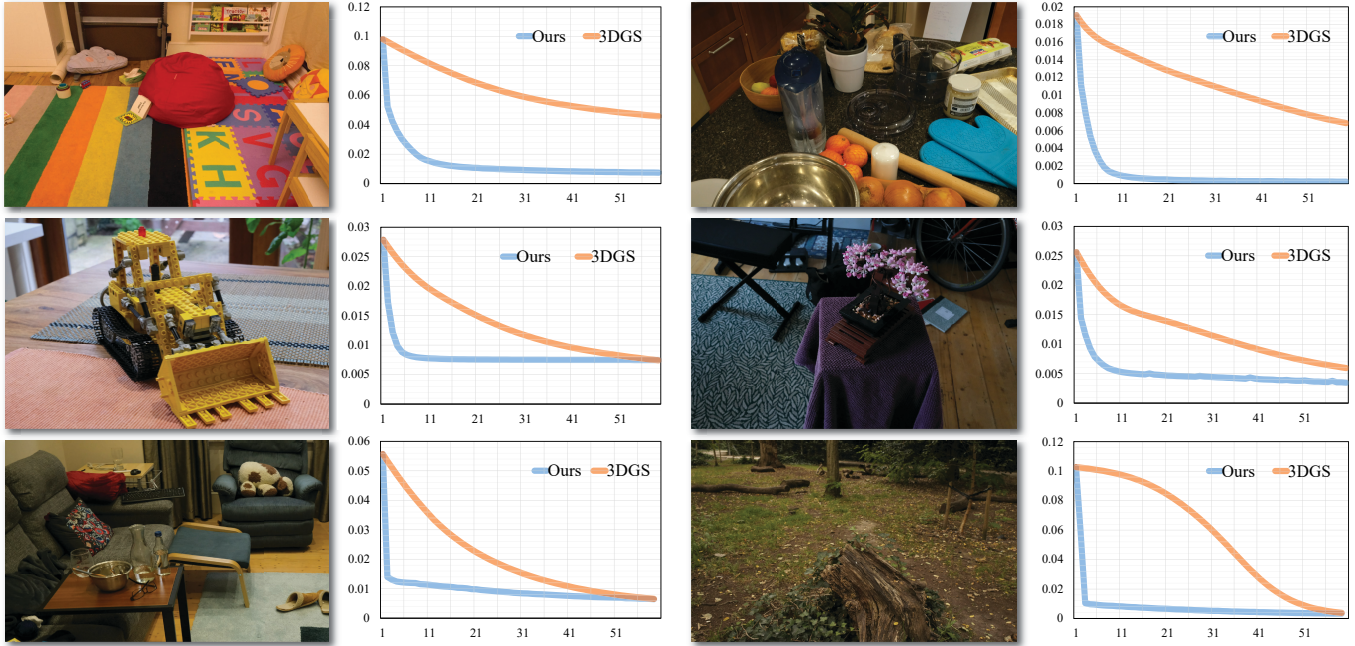
Fig. 4. **Convergence curves of our method and 3DGS**. We report a group of representative convergence plots using the proposed training method (local Newton) and vanilla 3DGS training (GD). The corresponding input training images are also attached next to the curves. We observe a strong second-order convergence using our method compared with gradient-based training in almost all the scenes.
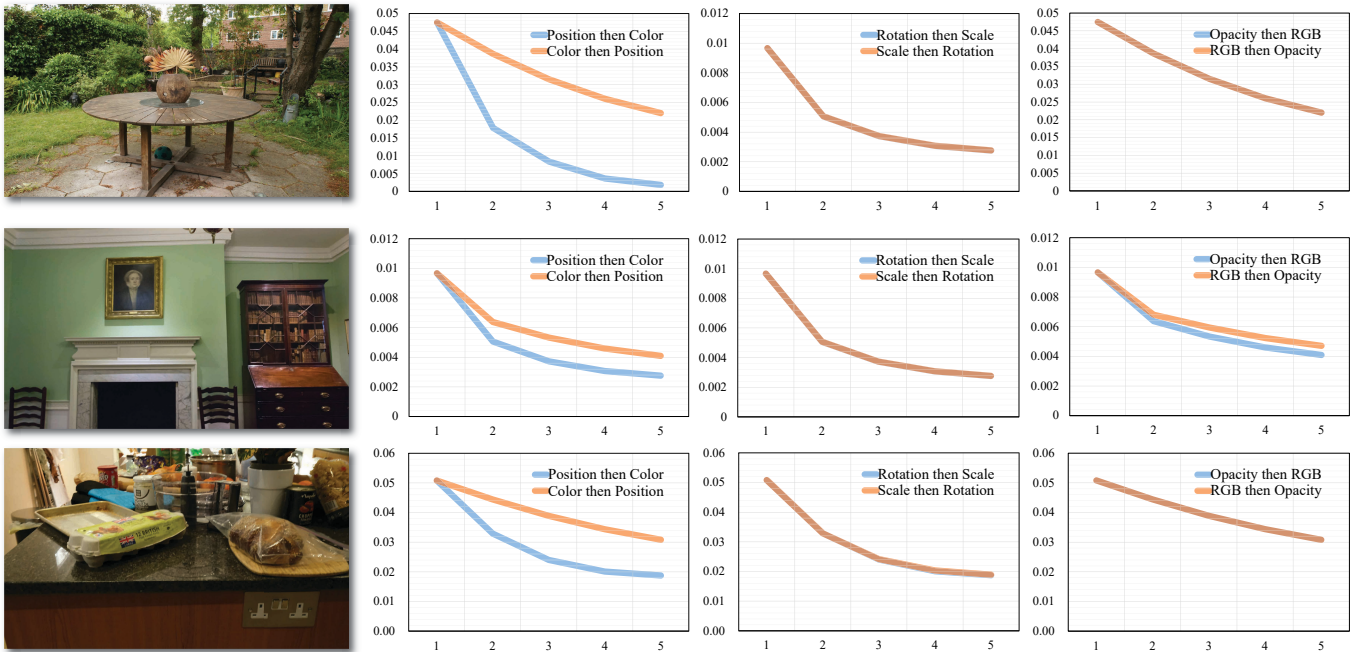


Fig. 5. **Training order of kernel attributes.** We always train the kernel position before the color information (RGB and opacity). Without a good positioning of the kernel, there is limited space for color-wise optimization. If one chooses to train the kernel color first, the training converges quadratically as well, but to a different local minimum. Meanwhile, the training order of rotation and scaling or opacity and RBG information is not important. The convergence curves are nearly identical as those attributes are de-coupled.

Table 1. **Benchmark statistics.** This table reports standard benchmarks using different 3DGS training algorithm, including vanilla 3DGS [Kerbl et al. 2023], AdR-Gaussian [Wang et al. 2024], EAGLES [Girish et al. 2025], 3DGS-LM [Höllein et al. 2024], and Taming 3DGS [Mallick et al. 2024]. The performance of our method can be further improved with fast rasterization techniques. The table only reports the performance of our method using the standard 3DGS rasterization pipeline.

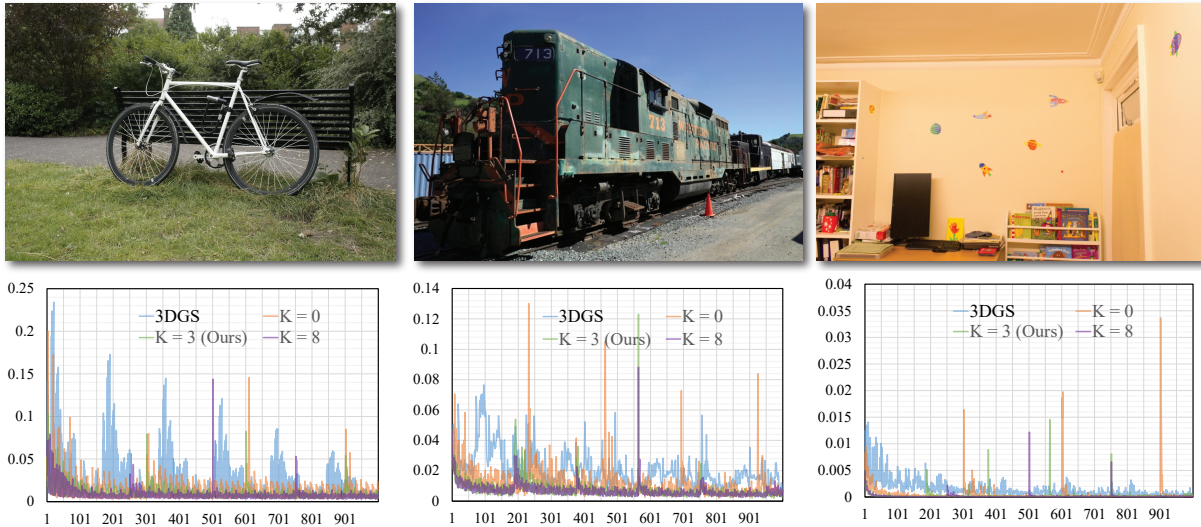| Training algorithm | Mip NeRF-360 | | | | Tanks & Temples | | | | Deep Blender | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | SSIM | PSNR | LPIPS | Training (s) | SSIM | PSNR | LPIPS | Training (s) | SSIM | PSNR | LPIPS | Training (s) |
| Vanilla 3DGS [Kerbl et al. 2023] | 0.871 | 29.18 | 0.183 | 1307 | 0.853 | 23.71 | 0.169 | 695 | 0.907 | 29.90 | 0.238 | 1210 |
| AdR-Gaussian [Wang et al. 2024] | 0.850 | 28.50 | 0.220 | 783 | 0.835 | 23.52 | 0.201 | 476 | 0.905 | 29.75 | 0.250 | 709 |
| EAGLES [Girish et al. 2025] | 0.810 | 28.64 | 0.192 | 2163 | 0.834 | 23.13 | 0.204 | 913 | 0.909 | 29.79 | 0.242 | 1560 |
| Taming 3DGS [Mallick et al. 2024] | 0.878 | 29.48 | 0.171 | 781 | 0.859 | 24.15 | 0.161 | 482 | 0.911 | 30.27 | 0.232 | 627 |
| 3DGS-LM [Höllein et al. 2024] | 0.813 | 27.39 | 0.221 | 972 | 0.845 | 23.72 | 0.182 | 663 | 0.903 | 29.72 | 0.247 | 951 |
| Our method | 0.876 | 29.42 | 0.168 | **256** (5.1×) | 0.871 | 24.43 | 0.154 | **131** (5.3×) | 0.927 | 30.43 | 0.228 | **189** (6.4×) |



Fig. 6. **Overshoot vs. $|\mathcal{N}_t|$.** The Hessian and gradient from secondary target $\mathcal{N}_t$ effectively avoids overshoot, making our training converge smoothly across different batches. The plots show that SGD-based 3DGS training overshoots — we can see loss spikes after the training switches to a new $I^t$. Without the secondary target i.e., $|\mathcal{N}_t| = 0$, , our method overshoots too, and sometimes even more severely than the vanilla 3DGS, due to its secondary-order nature. Fortunately, sampling at the secondary target of KNNs mitigates this issue. When $K = 8$, we barely see oscillations of the training loss. In our implementation, we use $K = 3$, which balances the training efficiency and convergence.

Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol (Eds.). Springer Nature Switzerland, Cham, 165–181.

Guofeng Feng, Siyan Chen, Rong Fu, Zimu Liao, Yi Wang, Tao Liu, Zhilin Pei, Hengjie Li, Xingcheng Zhang, and Bo Dai. 2024. FlashGS: Efficient 3D Gaussian Splatting for Large-scale and High-resolution Rendering. arXiv:2408.07967 [cs.CV]

Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. 2022. Plenoxels: Radiance Fields without Neural Networks. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 5491–5500.

Sharath Girish, Kamal Gupta, and Abhinav Shrivastava. 2025. EAGLES: Efficient Accelerated 3D Gaussians with Lightweight EncodingS. In *Computer Vision – ECCV 2024*, Aleš Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol (Eds.). Springer Nature Switzerland, Cham, 54–71.

Abdullah Hamdi, Luke Melas-Kyriazi, Jinjie Mai, Guocheng Qian, Ruoshi Liu, Carl Vondrick, Bernard Ghanem, and Andrea Vedaldi. 2024. GES: Generalized Exponential Splatting for Efficient Radiance Field Rendering. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 19812–19822.

Alex Hanson, Allen Tu, Geng Lin, Vasu Singla, Matthias Zwicker, and Tom Goldstein. 2024. Speedy-Splat: Fast 3D Gaussian Splatting with Sparse Pixels and Sparse Primitives. *arXiv preprint arXiv:2412.00578* (2024).

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.

Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul Debevec. 2021. Baking Neural Radiance Fields for Real-Time View Synthesis. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. 5855–5864.

Lukas Höllein, Aljaž Božič, Michael Zollhöfer, and Matthias Nießner. 2024. 3DGS-LM: Faster Gaussian-Splatting Optimization with Levenberg-Marquardt. arXiv:2409.12892 [cs.CV]

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.

Bernhard Kerbl, Andreas Meuleman, Georgios Kopanas, Michael Wimmer, Alexandre Lanvin, and George Drettakis. 2024. A hierarchical 3d gaussian representation for real-time rendering of very large datasets. *ACM Transactions on Graphics (TOG)* 43, 4 (2024), 1–15.

Shakiba Kheradmand, Daniel Rebain, Gopal Sharma, Weiwei Sun, Jeff Tseng, Hossam Isack, Abhishek Kar, Andrea Tagliasacchi, and Kwang Moo Yi. 2024. 3D Gaussian Splatting as Markov Chain Monte Carlo. arXiv:2404.09591 [cs.CV]

Diederik P Kingma. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

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.

Georgios Kopanas, Julien Philip, Thomas Leimkühler, and George Drettakis. 2021. Point-Based Neural Rendering with Per-View Optimization. *Computer Graphics Forum* 40, 4 (2021), 29–43.

Christoph Lassner and Michael Zollhöfer. 2021. Pulsar: Efficient Sphere-based Neural Rendering. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 1440–1449.

Joo Chan Lee, Daniel Rho, Xiangyu Sun, Jong Hwan Ko, and Eunbyung Park. 2024. Compact 3D Gaussian Representation for Radiance Field. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 21719–21728.

Sylvain Lefebvre and Hugues Hoppe. 2006. Perfect spatial hashing. *ACM Transactions on Graphics (TOG)* 25, 3 (2006), 579–588.

Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. 2019. Neural volumes: learning dynamic renderable volumes from images. *ACM Trans. Graph.* 38, 4, Article 65 (July 2019), 14 pages.

Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. 2021. Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–13.

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 *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 20654–20664.

Saswat Subhajyoti Mallick, Rahul Goel, Bernhard Kerbl, Markus Steinberger, Francisco Vicente Carrasco, and Fernando De La Torre. 2024. Taming 3DGS: High-Quality Radiance Fields with Limited Resources. In *SIGGRAPH Asia 2024 Conference Papers (SA '24)*. Association for Computing Machinery, New York, NY, USA, Article 2, 11 pages.

B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi, and R Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*.

Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant neural graphics primitives with a multiresolution hash encoding. 41, 4, Article 102 (July 2022), 15 pages.

Jorge Nocedal and Stephen J Wright. 1999. *Numerical optimization*. Springer.

Florian A Potra and Stephen J Wright. 2000. Interior-point methods. *Journal of computational and applied mathematics* 124, 1-2 (2000), 281–302.

Ning Qian. 1999. On the momentum term in gradient descent learning algorithms. *Neural networks* 12, 1 (1999), 145–151.

Ananth Ranganathan. 2004. The levenberg-marquardt algorithm. *Tutorial on LM algorithm* 11, 1 (2004), 101–110.

Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. 2021. KiloNeRF: Speeding up Neural Radiance Fields with Thousands of Tiny MLPs. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. 14315–14325.

Samuel Rota Bulò, Lorenzo Porzi, and Peter Kontschieder. 2025. Revising Densification in Gaussian Splatting. In *Computer Vision – ECCV 2024*, Aleš Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol (Eds.). Springer Nature Switzerland, Cham, 347–362.

Noah Snavely, Steven M Seitz, and Richard Szeliski. 2006. Photo tourism: exploring photo collections in 3D. In *ACM siggraph 2006 papers*. 835–846.

Kaiwen Song, Xiaoyi Zeng, Chenqu Ren, and Juyong Zhang. 2025. City-on-web: Real-time neural rendering of large-scale scenes on the web. In *European Conference on Computer Vision*. Springer, 385–402.

Cheng Sun, Min Sun, and Hwann-Tzong Chen. 2022. Direct Voxel Grid Optimization: Super-fast Convergence for Radiance Fields Reconstruction. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 5449–5459.

Justus Thies, Michael Zollhöfer, and Matthias Nießner. 2019. Deferred neural rendering: Image synthesis using neural textures. *Acm Transactions on Graphics (TOG)* 38, 4 (2019), 1–12.

Xinzhe Wang, Ran Yi, and Lizhuang Ma. 2024. AdR-Gaussian: Accelerating Gaussian Splatting with Adaptive Radius. In *SIGGRAPH Asia 2024 Conference Papers (SA '24)*. Association for Computing Machinery, New York, NY, USA, Article 73, 10 pages.

Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing* 13, 4 (2004), 600–612.

Vickie Ye, Ruilong Li, Justin Kerr, Matias Turkulainen, Brent Yi, Zhuoyang Pan, Otto Seiskari, Jianbo Ye, Jeffrey Hu, Matthew Tancik, and Angjoo Kanazawa. 2024. gsplat: An Open-Source Library for Gaussian Splatting. arXiv:2409.06765 [cs.CV]

Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. 2021. PlenOctrees for Real-time Rendering of Neural Radiance Fields . In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE Computer Society, Los Alamitos, CA, USA, 5732–5741.

Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. 2024. Mip-Splatting: Alias-Free 3D Gaussian Splatting. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 19447–19456.

Qiang Zhang, Seung-Hwan Baek, Szymon Rusinkiewicz, and Felix Heide. 2022. Differentiable Point-Based Radiance Fields for Efficient View Synthesis. In *SIGGRAPH Asia 2022 Conference Papers* (Daegu, Republic of Korea,). ACM, New York, NY, USA, Article 7.

Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. 2018. Stereo magnification: learning view synthesis using multiplane images. *ACM Trans. Graph.* 37, 4, Article 65 (July 2018), 12 pages.

## A DERIVATIVES OF INTERMEDIATE VARIABLES

We give the detailed formulation of the first and second derivatives of intermediate variables w.r.t. the kernel center $\boldsymbol{p}_k$.

### A.1 Projection function

$\boldsymbol{\pi}_k(\boldsymbol{p}_k)$ is the projection function, which converts the kernel center from 3D space to its 2D normalized device coordinate. Its first and second derivative w.r.t. the kernel center is:

$$\frac{\partial \boldsymbol{\pi}_k}{\partial \boldsymbol{p}_k} = \begin{bmatrix} \frac{W_I^t}{2}\left(\frac{1}{h_w}(PW)_0 - \frac{h_x}{h_w}(PW)_3\right)^\top \\ \frac{H_I^t}{2}\left(\frac{1}{h_w}(PW)_1 - \frac{h_y}{h_w}(PW)_3\right)^\top \end{bmatrix}, \tag{31}$$

$$\frac{\partial^2 \boldsymbol{\pi}_k}{\partial \boldsymbol{p}_k^2} = \begin{bmatrix} W_I^t\left(\frac{h_x}{h_w^2}(PW)_3^\top(PW)_3 - \frac{1}{h_w^2}(PW)_3^\top(PW)_0\right) \\ H_I^t\left(\frac{h_y}{h_w^2}(PW)_3^\top(PW)_3 - \frac{1}{h_w^2}(PW)_3^\top(PW)_1\right) \end{bmatrix}. \tag{32}$$

Here, $W_I^t$ and $H_I^t$ are the width and height of the input training $I^t$. $P$ and $W$ are $4 \times 4$ (homogeneous) projection matrix and viewing matrix.

$$\boldsymbol{h} = [h_x, h_y, h_z, h_w]^\top = PW[\boldsymbol{p}_k^\top, 1]^\top \in \mathbb{R}^4. \tag{33}$$

The notion $(PW)_i$ is a 3 dimensional *row vector* corresponding to the first three entries of $i$-th row of $PW$.

### A.2 SH color

$\tilde{\boldsymbol{c}}_k \in \mathbb{R}^3$ is the view-dependent color based on the SH coefficients. We give the derivatives for the red color component ($\tilde{c}_{k,R}$), and the formulation for the other two components is the same.

$$\frac{\partial \tilde{c}_{k,R}}{\partial \boldsymbol{p}_k} = \left(\boldsymbol{B}_{k,R} \odot \boldsymbol{c}_{k,R}\right)^\top \frac{\partial \Phi(\boldsymbol{r}_k)}{\partial \boldsymbol{r}_k} : \frac{\partial \Phi(\boldsymbol{r}_k)}{\partial \boldsymbol{p}_k},$$

$$\frac{\partial^2 \tilde{c}_{k,R}}{\partial \boldsymbol{p}_k^2} = \left(\frac{\partial \Phi(\boldsymbol{r}_k)}{\partial \boldsymbol{p}_k}\right)^\top : \left(\boldsymbol{B}_{k,R} \odot \boldsymbol{c}_{k,R}\right)^\top \frac{\partial^2 \Phi(\boldsymbol{r}_k)}{\partial \boldsymbol{r}_k^2} : \frac{\partial \Phi(\boldsymbol{r}_k)}{\partial \boldsymbol{p}_k}$$

$$+ \left(\boldsymbol{B}_{k,R} \odot \boldsymbol{c}_{k,R}\right)^\top \frac{\partial \Phi(\boldsymbol{r}_k)}{\partial \boldsymbol{r}_k} : \frac{\partial^2 \Phi(\boldsymbol{r}_k)}{\partial \boldsymbol{p}_k^2}. \tag{34}$$

$\boldsymbol{B}_k$ is the SH bases. $\odot$ represents the Hadamard product. $\boldsymbol{r}_k$ is the unit vector from the camera to $\boldsymbol{p}_k$. $\Phi(\boldsymbol{r}_k)$ is a set of vectors:

$$\Phi_{d=0} = 1,$$
$$\Phi_{d=1} = \left[\Phi_{d=0}, -r_{k,y}, r_{k,z}, -r_{k,x}\right]^\top,$$
$$\Phi_{d=2} = \left[\Phi_{d=0}, \Phi_{d=1}, r_{k,x}r_{k,y}, r_{k,y}r_{k,z}, 2r_{k,z}^2 - r_{k,x}^2 - r_{k,y}^2,\right.$$
$$\left. r_{k,x}r_{k,z}, r_{k,x}^2 - r_{k,y}^2\right]^\top,$$
$$\Phi_{d=3} = \left[\Phi_{d=0}, \Phi_{d=1}, \Phi_{d=2}, r_{k,y}(3r_{k,x}^2 - r_{k,y}^2), r_{k,x}r_{k,y}r_{k,z},\right.$$
$$r_{k,y}(4r_{k,z}^2 - r_{k,x}^2 - r_{k,y}^2), r_{k,z}(2r_{k,z}^2 - 3r_{k,x}^2 - 3r_{k,y}^2),$$
$$\left. r_{k,x}(4r_{k,z}^2 - r_{k,x}^2 - r_{k,y}^2), r_{k,z}(r_{k,x}^2 - r_{k,y}^2), r_{k,x}(r_{k,x}^2 - 3r_{k,y}^2)\right]^\top.$$

## A.3 Projected covariance matrix

Let $\boldsymbol{x} = [m, n]^\top$ be the image space coordinate of the pixel. $G_k = \exp^{-\frac{1}{2}(\boldsymbol{\pi}_k - \boldsymbol{x})^\top \Sigma^{-1} (\boldsymbol{\pi}_k - \boldsymbol{x})}$, the derivatives $\frac{\partial G_k}{\partial \Sigma_k}$, $\frac{\partial^2 G_k}{\partial \Sigma_k^2}$ can be written as follows:

$$\frac{\partial G_k}{\partial \Sigma_k} = -\frac{1}{2} G_k \left( (\boldsymbol{\pi}_k - \boldsymbol{x})^\top \frac{\partial \Sigma_k^{-1}}{\partial \Sigma_k} : (\boldsymbol{\pi}_k - \boldsymbol{x}) \right),$$

$$\frac{\partial^2 G_k}{\partial \Sigma_k^2} = \frac{G_k^2}{4} \left( (\boldsymbol{\pi}_k - \boldsymbol{x})^\top \frac{\partial \Sigma_k^{-1}}{\partial \Sigma_k} : (\boldsymbol{\pi}_k - \boldsymbol{x}) \right)^2$$

$$- \frac{G_k}{2} \left( (\boldsymbol{\pi}_k - \boldsymbol{x})^\top : \frac{\partial \Sigma_k^{-1}}{\partial \Sigma_k} : (\boldsymbol{\pi}_k - \boldsymbol{x}) \right). \tag{35}$$

$\frac{\partial \Sigma_k^{-1}}{\partial \Sigma_k}$ is a fourth-order tensor, and $\frac{\partial^2 \Sigma_k^{-1}}{\partial \Sigma_k^2}$ is a sixth-order tensor. The element-wise expressions can be written as:

$$\frac{\partial \left[ \Sigma_k^{-1} \right]_{i,j}}{\partial \left[ \Sigma_k \right]_{p,l}} = - \left[ \Sigma_k^{-1} \right]_{i,p} \left[ \Sigma_k^{-1} \right]_{l,j},$$

$$\frac{\partial \left[ \Sigma_k^{-1} \right]_{i,j}}{\partial \left[ \Sigma_k \right]_{p,l} \partial \left[ \Sigma_k \right]_{g,h}} = \left[ \Sigma_k^{-1} \right]_{i,g} \left[ \Sigma_k^{-1} \right]_{h,p} \left[ \Sigma_k^{-1} \right]_{l,j}$$

$$+ \left[ \Sigma_k^{-1} \right]_{i,p} \left[ \Sigma_k^{-1} \right]_{l,g} \left[ \Sigma_k^{-1} \right]_{h,j}.$$