# Tortho-Gaussian:
# Splatting True Digital Orthophoto Maps

Xin Wang, *Member, IEEE*, Wendi Zhang, Hong Xie, Haibin Ai, Qiangqiang Yuan, *Member, IEEE* and Zongqian Zhan, *Member, IEEE*

*Abstract*—True Digital Orthophoto Maps (TDOMs) are essential products for digital twins and Geographic Information Systems (GIS). Traditionally, TDOM generation involves a complex set of traditional photogrammetric process, which may deteriorate due to various challenges, including inaccurate Digital Surface Model (DSM), degenerated occlusion detections, and visual artifacts in weak texture regions and reflective surfaces, etc. To address these challenges, we introduce TOrtho-Gaussian, a novel method inspired by 3D Gaussian Splatting (3DGS) that generates TDOMs through orthogonal splatting of optimized anisotropic Gaussian kernel. More specifically, we first simplify the orthophoto generation by orthographically splatting the Gaussian kernels onto 2D image planes, formulating a geometrically elegant solution that avoids the need for explicit DSM and occlusion detection. Second, to produce TDOM of large-scale area, a divide-and-conquer strategy is adopted to optimize memory usage and time efficiency of training and rendering for 3DGS. Lastly, we design a fully anisotropic Gaussian kernel that adapts to the varying characteristics of different regions, particularly improving the rendering quality of reflective surfaces and slender structures. Extensive experimental evaluations demonstrate that our method outperforms existing commercial software in several aspects, including the accuracy of building boundaries, the visual quality of low-texture regions and building facades. These results underscore the potential of our approach for large-scale urban scene reconstruction, offering a robust alternative for enhancing TDOM quality and scalability. Project Web: https://gwen233666.github.io/Ortho-Gaussian/

*Index Terms*—3D Gaussian Splatting (3DGS), True Digital Orthophoto Maps, Occlusion Detection, Fully Anisotropic Gaussian Kernel.

## I. INTRODUCTION

**D**IGITAL Orthophoto Maps (DOM) not only encompass rich texture information but also exhibit the geometric properties inherent to maps, making them highly applicable in various fields[1], such as urban planning and cultural heritage preservation. Traditional DOMs generation takes oriented images and Digital Elevation Models (DEM) as input and employs digital differential rectification techniques. The generated result is a orthogonal nadir view of the surface that

Xin Wang and Wendi Zhang are co-first authors with equal contribution and importance. This work was supported by the National Natural Science Foundation of China (No. 42301507) and Natural Science Foundation of Hubei Province, China (No. 2022CFB727). (*Corresponding author: Zongqian Zhan*)

Xin Wang, Wendi Zhang, Hong Xie, Qiangqiang Yuan and Zongqian Zhan are with the School of Geodesy and Geomatics, Wuhan University, Wuhan, 430079, China PR. (xwang@sgg.whu.edu.cn, 2023202140056@whu.edu.cn, hxie@sgg.whu.edu.cn, yqiang86@gmail.com and zqzhan@sgg.whu.edu.cn).

Haibin Ai is with the Chinese Academy of Surveying & Mapping, Beijing, 100830, China PR. (aihb@casm.ac.cn)

Manuscript received April 19, 2005; revised August 26, 2015.

eliminates pin-hole projection distortions caused by terrain fluctuations and oblique photography, ensuring that the corresponding measuring scale is uniform throughout [2]. DOM exhibits artifacts and incorrect geometries that are resulted from occlusion by the façades of man-made architectures, as Fig. 1 shows. Therefore, the TDOM (True Digital Orthophoto Maps) is more extensively used, which takes DSM into account and perform visibility check of mesh triangles for detecting occlusion [3–5], and the obtained maps is less stemmed from the facades of buildings.

Over the past few decades, numerous traditional methods for True Digital Orthophoto Map (TDOM) generation have been developed to address occlusion detection. One of the most widely used techniques is Z-buffering [6–15], which records the distances between the perspective center and object points corresponding to image pixels, thereby determining visibility by selecting the nearest points. Habib et al. [13] introduced the Angle-Based method, employing adaptive radial and spiral sweeps to analyze the angles between lines connecting the perspective center and Digital Surface Model (DSM) cells. They also proposed a Height-Based method, which assesses ray height against ground points along the search path, identifying occlusions when any ground point exceeds the ray height. Kuzmin et al. [16] presented a Polygon-Based method that projects building polygons from the Digital Building Model (DBM) onto images, thereby enhancing orthographic image selection in sheltered areas. Wang et al. [17] applied a global variational model and texture matching techniques to fill in incomplete image data, thereby improving visual interpretation. Zhou et al. [18] developed a model linking ghost images to occlusions and employed seed growth algorithms for ghost detection.

In last few years, to the best of our knowledge, only a limited learning-based methods in scope have been proposed to enhance TDOM generation. Ebrahimikia and Hosseininaveh [19] employ a learning-based approach to detect building edges in images and estimate corresponding 3D edge points to modify the DSM. They subsequently propose urban-SnowflakeNet for completing point cloud of edge regions [20]. In addition, some works [21, 22] leverage the unoccluded properties of LiDAR intensity data within projection geometry to directly train a GAN (Generative Adversarial Network) model for generating TDOM. However, the first two methods suffer from limited generalization capabilities, while the latter heavily relies on intensity of LiDAR data. Furthermore, the emergence of neural radiance fields (NeRF) has revolutionized differentiable rendering by introducing an implicit 3D scene
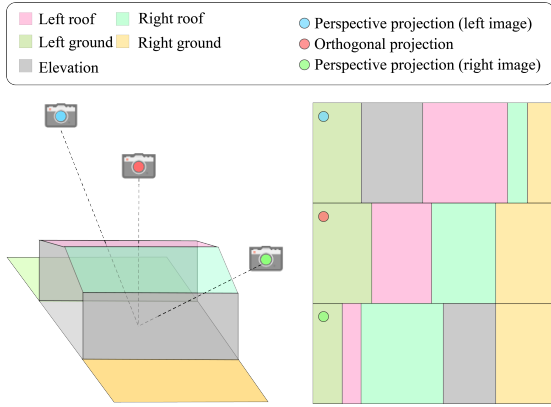
Fig. 1. A simplified model of building. It includes roofs and ground marked with different colors. The blue and green cameras represent perspective projection, with their imaging containing the gray building facades. The red camera represents orthographic projection.

representation through Multi-layer Perceptrons (MLPs) [23–30], presenting an alternative approach for TDOM generation. Chen et al. [31] introduce the Ortho-Nerf, which utilizes Plenoxels [32] (an accelerated version of NeRF), and a true-ortho-volume rendering strategy is employed for generating TDOM. Qu et al. [31] take satellite images as input and apply the RPC (Rational polynomial coefficients) instead of Structure from Motion (SfM) results to train a NeRF-based model to generate TDOM from satellite images using true-ortho-volume rendering.

Notwithstanding some advantages of novel view rendering, NeRF and its variants [32–35] are limited by the time-consuming training process and the high demands for real-time rendering. However, since the last year, 3D Gaussian Splatting (3DGS) has emerged as a promising alternative, offering efficient training and real-time rendering capabilities exceeding 100 frames per second (FPS), thereby becoming a research hotspot. As Fig. 2 illustrates, 3DGS employs a set of Gaussian Kernels (or ellipsoids) to explicitly represent 3D scene information, which can be reprojected onto a 2D image plane (known as "splatting"). The inherent rasterizer, implemented in CUDA, utilizes parallel-processing tiles to accelerate training and rendering. Consequently, it makes sense to apply 3DGS in the production of TDOMs. However, two significant challenges must be addressed: first, the scalability for large scenes. The number of 3D Gaussian kernels is bounded by the video memory, for example, a 24GB GPU can be used to optimize around 10 million Gaussians, while the small Garden scene of less than 100m2 in the Mip-NeRF360 [36] already needs about 5.8 million 3D Gaussians for rendering. Thus, scalability must be resolved to generate TDOM for large area; second, achieving high fidelity across various terrestrial scenarios. 3DGS often experiences blurring and aliasing when handling strong reflections and slender structures in aerial imagery, such as lakes or power lines.

In this work, we explore the potential of 3D Gaussian Splatting (3DGS) for generating True Digital Orthophoto Maps (TDOMs) through a novel method termed TOrtho-Gaussian. As Fig. 3 shows, we propose the first orthogonal splatting tech-

nique for rendering scale-uniform images, aka TDOMs, diverging from conventional perspective splatting. Additionally, inspired by VastGaussian [37], we adopt a divide-and-conquer strategy for large-scale scene 3DGS optimization, enabling the rendering of true orthophotos via orthogonal splatting. To achieve a high-fidelity TDOM for challenging scenarios, we introduce a plug-and-play Fully Anisotropic Gaussian Kernel (FAGK), which integrates transparency, scaling, and rotation parameters into spherical harmonic representations. This approach employs dense sampling to significantly enhance the depiction of highly reflective surfaces and slender structures. Our main contributions are as threefold:

1) The first attempt to generate TDOM via 3DGS, instead of the perspective splatting in the vanilla 3DGS, presenting a novel orthogonal splatting method which can geometrically elegant bypass the requirement of explicit DSM and occlusion detection.

2) The application of a divide-and-conquer strategy to extend scalability for generating TDOMs in large-scale scenes, which can improve both the time efficiency pf 3DGS optimization and the video memory usage.

3) The enhancement of 3D Gaussian representation by Fully Anisotropic Gaussian Kernel, incorporating transparency, scaling, and rotation parameters into spherical harmonic representations.

## II. RELATED WORKS

In this section, some previous relevant works are reviewed including the key challenge for generating TDOM – occlusion detection, learning-based and differentiable Rendering-Based TDOM generation methods.

### A. Occlusion Detection-Based TDOM

To date, ample endeavors have been dedicated to TDOM generation using traditional methods. However, many of these works continue to experience challenges such as misalignment, ghosting, and repeated mapping [9, 21–23, 38, 39]. These issues often stem from inaccuracy in DSM, but they are primarily exacerbated by insufficient occlusion detection during the digital differential rectification process. To address occlusion detection, various strategies have been developed, including those Z-buffering-based [6–15, 40], angles [13, 41], heights [9–11, 42, 43], vector polygons [9, 16, 38, 44–47], texture synthesis [17, 18], and object-oriented methods [18–20, 48].

The Z-buffer method is one of the earliest techniques introduced for generating orthophoto corrections based on Digital Models (DM) and Digital Building Models (DBM). The basic insight is that two orthophotos are first independently generated and subsequently merged into a single image. The visibility of 3D points is determined by resolving depth conflicts among points illuminated by the same light source, using a perspective projection center [4, 24, 48]. This method effectively manages competition among multiple object points along the same projection ray, designating the closest point as visible while marking the others as occluded. Numerous variants of the Z-buffer method have been proposed to enhance
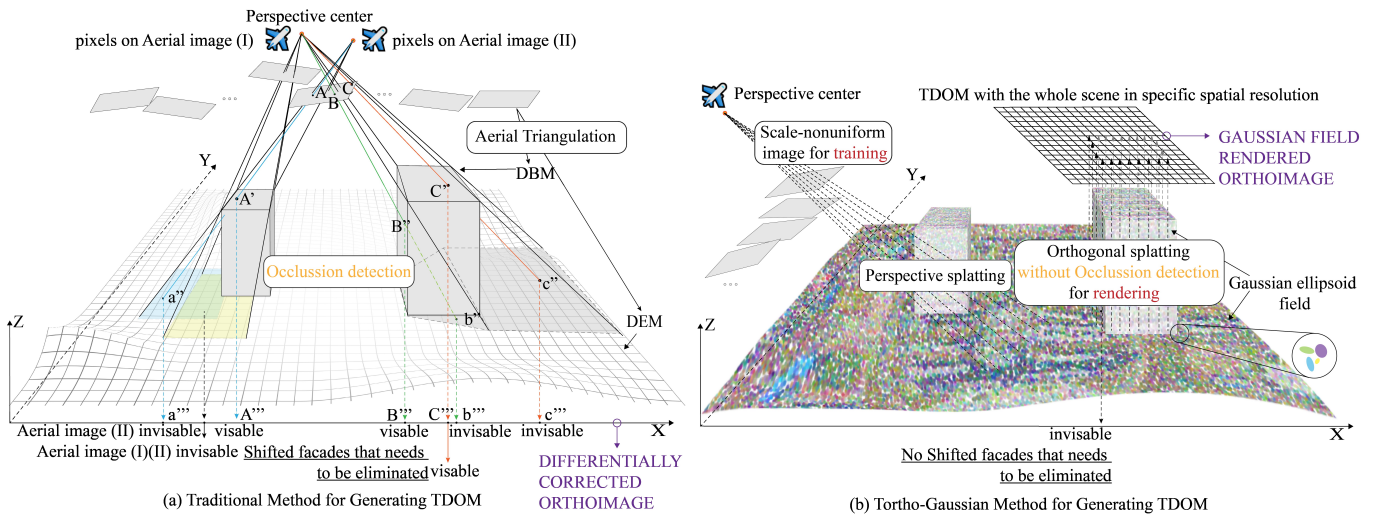
Fig. 2. Toy examples for generating TDOM. (a) The traditional photogrammetric TDOM generation from multi-view aerial images, incorporating an occlusion detection. Two images I and II are captured from two perceptive centers, with projection lines illustrating occlusion relationships from various viewpoints. In this solution, the DBM (Digital Building Model) and DEM (Digital Elevation Model) are required as input. These assist in occlusion detection, ensuring that only visible building surfaces are displayed, while eliminating façade shifts caused by viewpoint changes. (b) The proposed Tortho-Gaussian, generating a seamless, complete TDOM without the need for mosaicking. Scale-nonuniform and original images are used for 3DGS training. Then, an orthographic projection of the building cluster is performed along the z-axis, bypassing occlusion detection. The Gaussian ellipsoid field is rendered at a selected spatial resolution to produce a true orthophoto of the entire scene. This approach eliminates the need for post-processing steps such as image mosaicking and radiometric/color corrections, offering a streamlined and efficient solution.

its performance [12–18, 21, 22]. However, the Z-buffer method is sensitive to the ground sampling distance (GSD) of DSMs and is prone to the M-Portion problem when reconstructing elongated linear structures [4]. Consequently, many studies [7–15, 40] have focused on optimizing this method to address occlusion and pseudo-visibility issues.

To further mitigate artifacts, false occlusions, and pseudo-visibility, angle-based methodologies assess occlusion by evaluating the angular relationships between the camera rays of occluded and non-occluded regions in relation to vertical lines [13]. For example, Sheng et al. [41] develop a method for generating angle-based orthophotos of forest scenes using a Canopy Surface Model (CSM), Habib et al. [13] introduce an occlusion detection technique that calculates off-nadir angles between the lines of the perspective centers and the DSM pixels. Height-based methods have also been proposed, such as Habib et al.'s [10] height-based ranking approach and Bang & Kim's [15] height-based ray tracing method. These techniques detect occlusion by utilizing elevation information, identifying occluded regions by comparing building heights with the height of light rays along radial and helical paths [13].

Vector-based methods, such as the one proposed by Sheng et al. [38] , aim to eliminate artifacts generated by the Z-buffer technique via treating each pixel in the vector domain as a block rather than a single point. Additionally, Zhong et al. [44] introduced an occlusion detection method using polygon-based inversion imaging, leveraging the principle that building polygons do not occlude one another. In addition, texture synthesis methods have been utilized to compensate for missing image information, as demonstrated by Wang et al. [17], who employed a total variation model combined with texture matching. Wang et al. [25] further introduced

a line segment matching approach, where 2D line segments are extracted and matched, followed by 3D segment reconstruction, resulting in a highly accurate triangulated irregular network (TIN) model before TDOM generation. Li et al. [26] proposed a fusion algorithm based on the pulse-coupled neural network (PCNN) model. Lastly, Zhou et al. [18] developed an object-oriented model that addresses ghosting and occlusion artifacts through a seed growth method to detect occlusions in ghost images. Hu et al. [48] proposed a segmentation-based strategy for occlusion compensation, evaluating the simplicity of each segment's cost rather than relying on pixel-level quality assessments.

### B. Deep Learning-Based TDOM

In recent years, various efforts have focused on leveraging deep learning techniques to generate TDOM. Ebrahimikia and Hosseininaveh [49] propose a solution that improves the quality of DSM through pre-trained deep learning network. Their approach employs a 2D edge detector to identify building edges in imagery, which are then used to estimate 3D edges, thereby refining the DSM with explicit 3D edge points derived from the point cloud. Subsequently, Ebrahimikia et al. [50] later extend this approach by presenting the Urban-SnowflakeNet, aimed at completing building point clouds from photogrammetric processing. However, these approaches are limited to structural buildings, lacks generalizability, and still fails to address occlusion detection. Shin et al. [21, 22] enhance the quality of true orthophotos by employing the Pix2Pix model within a Generative Adversarial Network (GAN). Their approach primarily leverages LiDAR data to generate orthophotos that are free from projection geometry occlusions. Nevertheless, their methods heavily depend on the quality of the preprocessed LiDAR intensity.
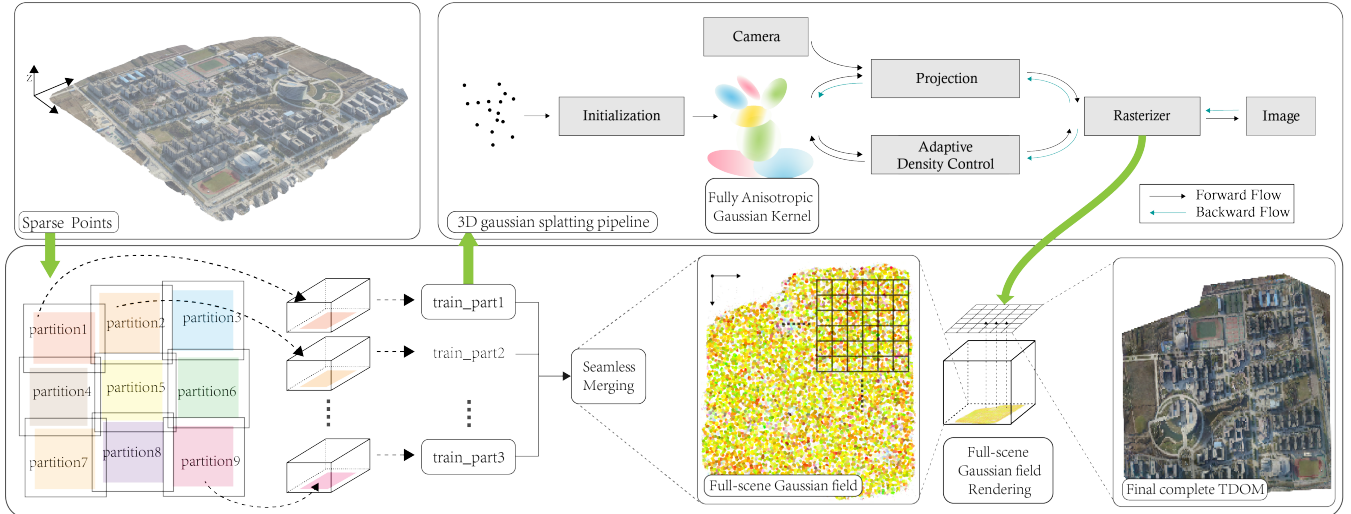
Fig. 3. Workflow of the proposed Tortho-Gaussian. First, it begins by aligning the sparse point cloud to the *x*- and *y*- axes. The scene is then partitioned into smaller regions using a divide-and-conquer strategy. For each partition, the original 3D Gaussian Splatting (3DGS) framework is employed to train the corresponding Gaussian field with the enhanced Fully Anisotropic Gaussian Kernel (FAGK). The trained Gaussian fields are seamlessly merged into a unified field to represent the entire scene, eliminating the need for TDOM tiling and subsequent color and brightness balancing steps. Next, the camera position is selected, and 3D Gaussian ellipsoids are splatted orthogonally, followed by pixel rasterization to compute the color for each pixel on the TDOM of the complete scene.

## C. Differentiable Rendering-Based TDOM

The rapid advancements in NeRF have profoundly transformed the 2D image reconstruction methods. However, to the best of our knowledge, it is two years later than the NeRF [51] work was proposed that implicit methods began to be leveraged for the generation of true digital orthophotos. Lv et al. [52] are the first to employ implicit neural representations for the rapid generation of digital orthophotos, exploring the potential of implicit reconstruction techniques utilizing the speed-optimized Instant NGP as a baseline. This approach holds the potential to fundamentally alter traditional orthophoto generation practices. Expanding on this groundwork, Chen et al. [53] employed the speed-optimized Plenoxels [32] during the rendering phase, using orthogonal projection to directly produce high-quality TDOM. This method effectively circumvents the intricate challenges associated with visibility analysis and texture compensation that are commonly encountered in conventional approaches, while also mitigating issues related to terrain edge distortion. Later, Qu et al. [31] further extend this idea by inputting satellite images for generating more large-arear TDOM.

Recently, 3DGS, as a differentiable rendering method that incorporates explicit spatial structural information, has started to supplant tasks accomplished by NeRF [54–57]. This shift is attributed to 3DGS's exceptional speed and high-detail novel view synthesis quality [58, 59]. Its flexible and adaptive 3D object representation overcomes the limitations of traditional volumetric rendering methods [56]. Therefore, this study seeks to explore the feasibility and potential of 3DGS for generating TDOM.

## III. METHODOLOGY

This section provides detailed explanations of our Tortho-Gaussian on generating TDOM, which contains five technical parts, i.e., preliminaries of 3GDS (III-A), orthographic splatting of 3DGS (III-B), TDOM generation based on orthographic splatting (III-C), fully anisotropic Gaussian Kernels (III-D) and the divide-and-conquer strategy (III-E) for large-scale area. More specifically, as Fig. 3 shows, after structure from motion, our Tortho-Gaussian first divide the whole scene into sub-regions, which are optimized by the 3DGS with an improved fully anisotropic Gaussian Kernel. After the optimization, we algin the 3DGS field such that the scene is parallel to the xy-plane. Then, the centroid of all camera are explored to set the spatial resolution, and performing orthogonal splatting along the z-axis to produce a true orthophoto of the complete scene, without the need for image stitching or color correction.

## A. Preliminaries

Following portion of [60], to make this paper more self-contained, we next outline some basics of 3D gaussian splatting. The 3DGS represents the scene using a series of Gaussian ellipsoid kernels that are closely aligned with the 3D scene's structures. Each Gaussian kernal is defined by a set of attributes, including the position (mean) $\mu$, anisotropic covariance $\Sigma$, opacity $\alpha$, and color c that is formulated via spherical harmonics. During the rendering stage, the position and covariance attributes of all Gaussians in the 3D scene are reprojected onto the image plane (namely Splatting), thereby forming 2D Gaussians. The reprojection of 3D Gaussians onto specific image tiles is determined based on the position and radius of these 2D Gaussians. The rendered image is then generated using a volumetric rendering technique combined with alpha blending. Ultimately, the Gaussian kernels are optimized based on the discrepancy between the rendered image and the input image. A Gaussian kernel $G_\Sigma(\mathbf{x})$, centered at $\mu$, with a

3D covariance matrix given by $\Sigma$ is represented via following formula.

$$G_\Sigma(\mathbf{x} - \mu) = \frac{1}{(2\pi)^{3/2} \mid \Sigma \mid^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)} \quad (1)$$

Where x and $\mu$ represent column vectors, specifically $[x, y, z]^T$ and $[\mu_x, \mu_v, \mu_z]^T$, respectively. The term $\Sigma$ denotes a symmetric $3 \times 3$ matrix. In order to maintain its positive semi-definite property, the covariance matrix is further parameterized by a scaling matrix S and a rotation matrix R:

$$\Sigma = RSS^T R^T \quad (2)$$

To perform the splatting, the 3D Gaussian kernel is projected onto a 2D Gaussian ellipse. The corresponding 2D covariance matrix $\Sigma'$ is calculated by the following formula:

$$\Sigma' = JW\Sigma W^T J^T \quad (3)$$

W denotes the viewpoint transformation, and $\Sigma$ represents the 3D covariance matrix, while J denotes the Jacobian matrix associated with the projective transformation within its affine approximation [60, 61].

For a specific pixel, the corresponding depths of all intersecting Gaussians can be derived with the viewing transformation W, which are employed to sort these intersected Gaussians $\mathcal{N}$. Subsequently, the final color of the pixel is calculated using alpha blending:

$$C = \sum_{n=1}^{|\mathcal{N}|} c_n \alpha'_n \prod_{j=1}^{n-1} \left(1 - \alpha'_j\right) \quad (4)$$

$c_n$ denotes the predicted color. The final splatting opacity $\alpha'_n$ is obtained by multiplying the predicted opacity $\alpha_n$ with the splatted 2D Gaussian, as defined below:

$$\alpha'_n = \alpha_n e^{-\frac{1}{2}(x'-\mu'_n)^T \Sigma'^{-1}_n (x'-\mu'_n)} \quad (5)$$

x' and $\mu'_n$ are coordinates defined in the 2D image plane. .

### B. Orthogonal Splatting of 3DGS

The 3D Gaussian Splatting technique offers significant potential for accurate scene representation and rendering, similar to the advancements in TDOM generation seen with NeRF [52, 53]. Based on the vanilla 3DGS that employs perspective splatting, we introduce orthogonal splatting which mainly includes the projection of the mean and variance of 3D Gaussian.

*1) Projection of 3D Gaussian - Mean:* The Gaussian mean represents the position coordinates of the Gaussian kernel, which is capable of undergoing a non-affine transformation through perspective projection. The corresponding transformation matrix for this projection is given by:

$$P = \begin{pmatrix} \frac{2z_n}{r-l} & 0 & \frac{r+l}{r-l} & 0 \\ 0 & \frac{2z_n}{t-b} & \frac{t+b}{t-b} & 0 \\ 0 & 0 & -\frac{z_f+z_n}{z_f-z_n} & -\frac{2z_f z_n}{z_f-z_{near}} \\ 0 & 0 & -1 & 0 \end{pmatrix} \quad (6)$$

where

$$r = \tan\left(\frac{\theta_x}{2}\right) \cdot z_n, \, l = -r \quad (7)$$

$$t = \tan\left(\frac{\theta_y}{2}\right) \cdot z_n, \, b = -t \quad (8)$$

$z_n$ and $z_f$ represent the near and far clipping planes, respectively, while $\theta_x$ and $\theta_y$ are the horizontal and vertical fields of view. $l$, $r$, $b$ and $t$ denote the left, right, top, and bottom boundaries of the viewing frustum. The Gaussian kernel is projected onto the 2D image plane through this perspective projection.

In order to make an orthographic projection for the mean value, equation (6) can be replaced by the following formula:

$$P_o = \begin{pmatrix} \frac{2}{r-l} & 0 & 0 & -\frac{r+l}{r-l} \\ 0 & \frac{2}{t-b} & 0 & -\frac{t+b}{t-b} \\ 0 & 0 & -\frac{2}{z_f-z_n} & -\frac{z_f+z_n}{z_f-z_n} \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (9)$$

Via equation (9), the center of the Gaussian sphere is orthographically splatted into a corresponding 2D Gaussian.

*2) Projection of 3D Gaussian Axes (Covariance):* [62] proved that covariance matrix can be represented using rotation and scaling matrices. In the context of projective transformations, the expected perspective projection matrix is given by formula (6). It is evident that this transformation is nonlinear and non-affine. To approximate the affine transformation of the Gaussian ellipsoid, a corresponding Jacobian matrix is employed for local linear approximation:

$$J = \begin{pmatrix} \frac{focal_x}{t_z} & 0 & -\frac{focal_x \cdot t_x}{t_z^2} \\ 0 & \frac{focal_y}{t_z} & -\frac{focal_y \cdot t_y}{t_z^2} \\ 0 & 0 & 0 \end{pmatrix} \quad (10)$$

where $focal_x$ and $focal_y$ represents the focal lengths of the camera along the x and y axes. $(t_x, t_y, t_z)$ is the coordinates of a 3D point in camera space and $(\nu_x, \nu_y, \nu_z, 1)$ is the corresponding homogeneous coordinate. Then, the Jacobian matrix $J_o$ corresponding to the orthogonal matrix $P_o$ can be derived by formula (11).

$$\nu' = P_o \begin{pmatrix} \nu_x \\ \nu_y \\ \nu_z \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{2}{r-l}\nu_x + \frac{r+l}{r-l} \\ \frac{2}{t-b}\nu_y + \frac{t+b}{t-b} \\ -\frac{2}{z_f-z_n}\nu_z + \frac{z_f+z_n}{z_f-z_n} \\ 1 \end{pmatrix} \quad (11)$$

Differentiating this equation, we can obtain:

$$J_o = \begin{pmatrix} \frac{\partial \nu'_x}{\partial \nu_x} & \frac{\partial \nu'_x}{\partial \nu_y} & \frac{\partial \nu'_x}{\partial \nu_z} \\ \frac{\partial \nu'_y}{\partial \nu_x} & \frac{\partial \nu'_y}{\partial \nu_y} & \frac{\partial \nu'_y}{\partial \nu_z} \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} \frac{2}{r-l} & 0 & 0 \\ 0 & \frac{2}{t-b} & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad (12)$$

To accurately orthogonally splat a 3D Gaussian onto a referenced 2D plane, both the Gaussian mean and variance have to undergo appropriate orthographic projection to ensure correct coverage of the corresponding tiles and pixels. Subsequently, this allows each pixel to be encoded as related 2D Gaussian in an orthographic manner, which are then rendered using $\alpha$-blending based on the depth order.

## C. TDOM generation based on orthographic splatting

To generate the TDOM of the whole scene, we have to perform an accurate orthographic projection for the entire scene, while avoiding the need for image mosaicking. First, the target spatial resolution $s_x$ and $s_y$ for each pixel on the TDOM should be set. During orthographic splatting, the grid points (as shown in Fig. 2 and 3), formed by the pre-set spatial resolution, are applied as referenced coordinates for rasterization. The 3D Gaussians that locate within one specific grid are employed to render the color information of corresponding TDOM pixel. The coordinates of each TDOM pixel can be defined as follows:

$$X = \left\{ \bar{X} + s_x \cdot (i - (W/2) + \delta_x) \mid i = 0, 1, \ldots, W \right\} \quad (13)$$

$$Y = \left\{ \bar{Y} + s_y \cdot (j - (H/2) + \delta_y) \mid j = 0, 1, \ldots, H \right\} \quad (14)$$

where, $W$ and $H$ is the width and Height of the TDOM which is up to the spatial resolution, $\delta_x$ and $\delta_y$ is constant value (equal to 1/2 spatial resolution) to make projection ray stay at the center of the pixel. For the position coordinates of $N$ cameras, denoted as $\{(x_i, y_i, z_i) \mid i = 1, 2, \ldots, N\}$, the center of the orthographic image, denoted as $(\bar{X}, \bar{Y})$ can be calculated using the centroid of the cameras' positions in the $x$- and $y$- directions:

$$\bar{X} = \left\{ \frac{1}{N} \sum_{i=1}^{N} x_i \right\} \quad (15)$$

$$\bar{Y} = \left\{ \frac{1}{N} \sum_{i=1}^{N} y_i \right\} \quad (16)$$

Finally, for each pixel, we orthographically splat the corresponding 3D Gaussians into 2D Gaussians, which are then $\alpha$-blended based on the depth sort for esimating the color. The complete TDOM is obtained by running all pixels as above.

## D. Fully Anisotropic Gaussian Kernel

Inspired by MLP-based Gaussian kernels [63], we develop a novel fully anisotropic Gaussian kernel that incorporates spherical harmonic (SH) coefficients for all corresponding properties, including color, opacity, rotation, and scaling. This representation enables the Gaussian kernel at a given position to exhibit distinct colors $c$, opacities $\alpha$, rotations $r$, and scales $s$ when observed from different directions, allowing our approach to achieve comparable performance to neural Gaussians [63].

Specifically, the direction between the Gaussian kernel and the camera center is first calculated resulting in direction-dependent Gaussian properties. The complexity of the spherical harmonics depends on the degree of the SH coefficients. In this work, similar to 3DGS, the coefficients up to the third order are applied, and upgrade the order every 1000 iterations. For each channel of Gaussian properties, the spherical harmonics are formulated as:

$$A_l^m(\theta, \varphi) = \begin{cases} \sqrt{2} K_l^m cos(m\varphi) P_l^m(cos\theta) & m > 0 \\ \sqrt{2} K_l^m sin(-m\varphi) P_l^{-m}(cos\theta) & m < 0 \\ K_l^0 P_l^0(cos\theta) & m = 0 \end{cases} \quad (17)$$

where

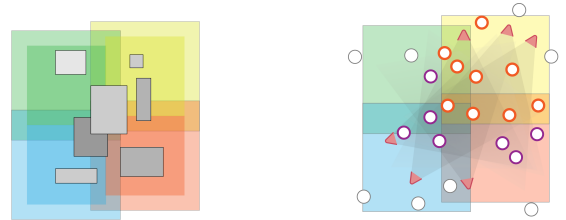$$P_n(x) = \frac{1}{2^n \cdot n!} \frac{d^n}{dx^n} \left[ (x^2 - 1)^n \right]$$

$$P_l^m = (-1)^m (1 - x^2)^{\frac{m}{2}} \frac{d^m}{dx^m} (P_l(x)) \quad (18)$$

$$K_l^m = \sqrt{\frac{(2l+1)(l - \mid m \mid)!}{4\pi(l + \mid m \mid)!}}$$

$K_l^m$ is the normalization constant, with $l$ as the degree and $m$ as the order of the spherical harmonics. $A_l^m(\theta, \phi)$ is the general form of the spherical harmonics. $P_l^m(x)$ is the associated Legendre polynomial. $\theta$ and $\varphi$ are the polar and azimuthal angles in spherical coordinates, respectively .

## E. Divide and Conquer strategy

The vanilla 3DGS uses explicit 3D Gaussian ellipsoids as primitives to represent a 3D scene. However, it is limited in handling only a finite number of 3D Gaussian ellipsoids, which is typically feasible for small-scale objects [31, 37, 64, 65]. As the scale of the scene increases, the demand for GPU memory (VRAM) grows substantially. For example, given a 24GB RTX 3090 GPU, it runs into memory shortage when the number of Gaussians exceeds approximately 10 million. Reconstructing large-scale scenes typically requires significantly higher number of Gaussians, further exacerbating this issue [65]. Additionally, the increasement of Gaussians also leads to a considerable slowdown in the depth sorting process prior to rendering.



(a) Scene Partitioning and Expansion    (b) Camera and Point Selection

Fig. 4. Divide-and-Conquer strategy. The dark rectangles are the initial divisions of the scene, the light-colored rectangles indicate the extended regions, and the orange and purple dots represent cameras selected by the local region and the external region, respectively.

To address this issue, analogous to [37], we adopted a divide-and-conquer strategy, introducing a progressive partition approach. Basically, the key insight is that the scene is divided into multiple overlapping rectangular regions for optimization with limited VRAM, reducing both training time and memory requirement. Fig. 4 illustrates the strategy we employ, and four steps are performed for the division:

*1) Scene Partitioning Based on Camera Positions:* In general, the input images are distributed across the scene to facilitate comprehensive reconstruction. Therefore, dividing the scene based on the distribution of camera centers makes sense. The camera centers are projected onto the ground plane. An $m \times n$ grid is adopted and ensures each cell contains a similar number of training images, thereby balancing the optimization procedure across different cells. As [37], the views $V$ on ground plane is divided into $m$  sections, each

(1) phantom3-npu
(2) phantom3-ieu
(3) phantom3-freeway
(4) phantom3-huangqi
(5) phantom3-hengdong
(6) phantom3-village
(7) phantom3-factory
(8) phantom3-centralPark
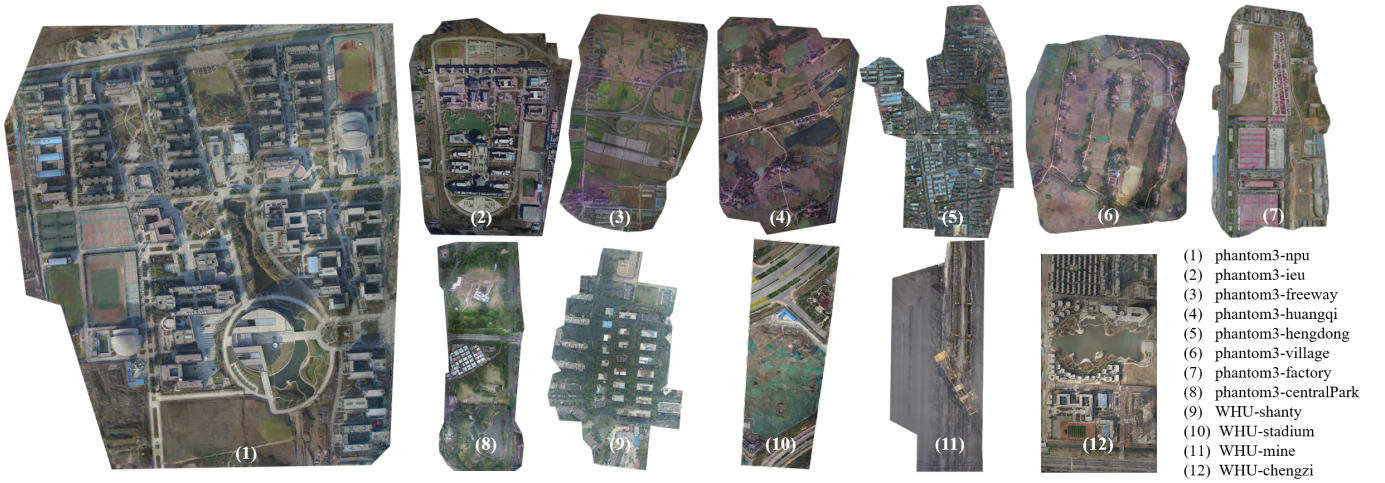(9) WHU-shanty
(10) WHU-stadium
(11) WHU-mine
(12) WHU-chengzi

Fig. 5. The complete TDOMs of the NPU-DroneMap and WHU dataset using our proposed Tortho-Gaussian.

containing approximately $|\mathbf{V}|/m$ views. These sections are further subdivided into $n$ cells, with around $|\mathbf{V}|/(m \times n)$ views per cell.

*2) Sub-division expansion:* To guarantee the consistency between adjacent subdivisions, it is necessary to expand the partitioned cells to some degrees. As Fig. 4 (a) illustrates that the light-colored rectangles are from the extended darker ones. For the $j^{th}$ cell, defined within a rectangle of size $\ell_i^h \times \ell_i^w$, we expand along the boundary by 20%, leading to a larger cell. This expanded boundary ensures more comprehensive coverage which contained extended subsets of camera set $\{\mathbf{V}_i\}_{i=1}^{m \times n}$ and point set $\{\mathbf{P}_i\}$ for each cell, allowing for improved continuity between adjacent cells,

*3) Camera Selection Based on Visibility:* Incorporating more images generally provides additional training samples, thereby enhancing the performance of 3DGS. For each divided cell, Fig. 4 (b) shows the selection of additional cameras based on visibility to enhance reconstruction fidelity. Cameras with visibility values exceeding a predefined threshold $t$ are chosen, where visibility is defined as the ratio of the projected area of the divided cell onto an image to the area of the entire image.

*4) Point Selection Based on Coverage:* After adding relevant cameras to the cell's camera set $\mathbf{V}_i$, all the points observed by these extra selected cameras are added to the point set $\mathbf{P}_i$. This step ensures a better initialization for the optimization of each cell. Proper initialization helps mitigate depth ambiguities, which could otherwise lead to incorrect 3D Gaussian distributions when fitting objects, particularly those located outside the cell. This method ensures accurate 3D Gaussian generation and prevents the formation of floating artifacts.

## IV. EXPERIMENTS

In this section, comprehensive experimental results are reported to demonstrate the efficacy of the proposed Tortho-Gaussian. Based on various datasets, both qualitative and quantitative evaluations are conducted via comparison with several state-of-the-art commercial software and ablation studies.

### A. Experimental settings

**Experiments protocols.** Overall, our experiments mainly contains three parts: first, Qualitative Performance. We compare our method with several state-of-the-art (SOTA) commercial software. Specifically, the quality of the generated TDOM is evaluated by analyzing building edges, facades, slender structures, and regions with weak textures; second, quantitative performance. To assess the precision of the generated TDOM, we first evaluate relative precision using Metashape and Pix4DMapper as benchmarks. Additionally, we provide an overlaid comparison with manually generated CAD maps to quantify the absolute mapping precision; Finally, ablation studies. We conduct extensive ablation studies to explore various aspects of our Tortho-Gaussian method, including spatial resolution, partitioning strategies, and the use of fully anisotropic Gaussian kernels. This comprehensive evaluation highlights the robustness and effectiveness of our approach.

**Experimental Data.** In our work, we first employ the NPU DroneMap dataset published by Bu et al. [66], it was captured images across different areas in China using a custom-built hexacopter equipped with a Phantom3 camera. Additionally, we generate a self-constructed dataset, WHU, including Shanty, Stadium, Mine and chengzi datasets using on Canon EOS 5D Mark III under different flight Height, among which the chengzi dataset contain a high-precision CAD map that is made by manual labeling. These experimental data consists of high-resolution aerial imagery captured in urban, rural, and mixed environments, encompassing various scene elements such as buildings, vegetation, roads, and water bodies. The overall views of these datasets are shown in Fig. 5, more detailed information can be found in Tab. I.

**Experimental Details**. Before 3D Gaussian optimization, the sparse point cloud of SfM was first preprocessed by Manhattan alignment, ensuring that the x and y axes of the point cloud were parallel to the boundary frame. Then,

TABLE I
DETAILS OF OUR EXPERIMENTAL DATASETS

| Dataset | Sequence | Location | H-Max (m) | Area (km²) | Imagenum | Imagesize (pix) |
|---------|----------|----------|-----------|-----------|----------|-----------------|
| NPU-DroneMap | phantom3-hengdong | Hengdong, Hunan | 358.00 | - | 221 | 1920*1080 |
| | phantom3-huangqi | Hengdong, Hunan | 222.30 | 1.313 | 393 | 1920*1080 |
| | phantom3-centralPark | Shenzhen, Guangdong | 161.80 | 0.606 | 835 | 1920*1080 |
| | phantom3-factory | - | 198.72 | 0.912 | 402 | 1920*1080 |
| | phantom3-freeway | - | 258.30 | 1.457 | 415 | 1920*1080 |
| | phantom3-village | Hengdong, Hunan | 160.60 | 0.932 | 406 | 1920*1080 |
| | phantom3-ieu | Zhenzhou, Henan | 282.30 | 1.524 | 467 | 1920*1080 |
| | phantom3-npu | Xi'an, Shaanxi | 254.50 | 1.598 | 457 | 1920*1080 |
| WHU | shanty | Wuhan, Hubei | 120.00 | 0.503 | 67 | 1600*900 |
| | stadium | Tianjin | 310.00 | - | 230 | 1600*900 |
| | mine | Rizhao, Shandong | 60.60 | - | 43 | 1600*900 |
| | chengzi | Suqian, Jiangsu | 167.35 | - | 819 | 1228*820 |

*'-' denotes missing information.*

the resulting transformation matrix was applied to each sub-regional training. During optimization, for every fixed number of images, one image was selected as a test image. All experiments were conducted on four NVIDIA GeForce RTX 4090 GPUs, with the number of iterations set to 30,000.

### B. Qualitative Evaluation

Four commercial software (ContextCapture [67], Metashape [68], Pix4DMapper [69], and Map2DFusion [66]), incorporated with traditional photogrammetric techniques relying on DSM, are compared on the NPU DroneMap and WHU dataset regarding the quality of TDOM. All these methods use the poses of ContextCapture before conducting the subsequent reconstruction tasks. Two indicators that can significantly reflect the quality of TDOM are visually investigated, i.e., building edges and facades.

*1) Building Edges:* In general, the edges of buildings on a satisfactory TDOM should follow the real geometry of building without irregular deformation, and the joints between buildings should precisely align. Fig. 6 and Fig. 7 shows the results of TDOM in a region densely packed buildings, and several different methods are compared. As it can be seen from Fig. 6 (e) and Fig. 7 (a), the reconstructed edges in Map2DFusion, ContextCapture, Metashape and Pixel4DMapper all exhibit varying degrees of distortion due to the insufficient digital differential correction. Even the Nerf-based method Ortho-NeRF [53] shows some artifacts that may be influenced by the inherent rendering performance of neural-based implicit 3D representation. On the other hand, our Tortho-Gaussian successfully generates high-fidelity building edges.

To further explore the quality of edge reconstruction, six buildings with prominent straight lines are selected from the NPU Dronemap and WHU datasets for extra assessment. In particular, we applied the Canny edge detection algorithm [70], using minimum and maximum thresholds, to extract line edge structures from the building outlines[1]. Based on the discrete points extracted on the generated TDOM, the lines are fitted using least-squares in an iterative manner. The points that stay within 2.5 time the standard deviation are considered as inliers and used to reflect quality of line edges obtained by various methods.

As shown in Fig. 8 and Fig. 9, comparing to the other methods, in general, our method can reduce the number of noise points in the linear building edges and achieve SOTA performance. Our Tortho-Gaussian is basically on par with other methods, wherein we yield the best results on some tests and others generate the best on other datasets.

These experimental results demonstrate that our approach can effectively suppresses jagged edges and curvature along building outlines.

*2) Building Facades:* For true orthophotos, the side facades of building should appear as just continuous boundaries, without breaks or jagged edges.Typically, the extent to which building facades can be fully resolved reflects the effectiveness of occlusion detection in the true orthophoto generation process. Our method effectively eliminates almost all facade occlusions. This can be clearly observed in Fig. 6 and Fig. 7. Fig. 6 (a), (b), (c) and (d) all exhibit prominent building facades, losing the characteristics of true orthographic photos. In Fig. 6(a,4), a blurred thick edge appears on the left facade, indicating significant distortion. Similar degradation can be also found in Fig. 7 (a),(b), (c). However, the Tortho-Gaussian method can successfully retain continuous linear building edges while completely eliminating the projections of building facades.

*3) Slender structures:* Reconstructing slender and thin structures, such as non-rigid tree canopies, trunks, and cables, is particularly challenging due to the inherent artifacts in novel view synthesis [71]. To address these issues, we employ the fully anisotropic Gaussian kernel functions as a super-resampling technique specifically designed to handle structures

---

[1]Empirically, we set the minimum threshold of the Canny algorithm as 50% of the maximum threshold to detect primary edge features with significant pixel gradients. The dimension of the Sobel operator, or aperture size, was set to 3
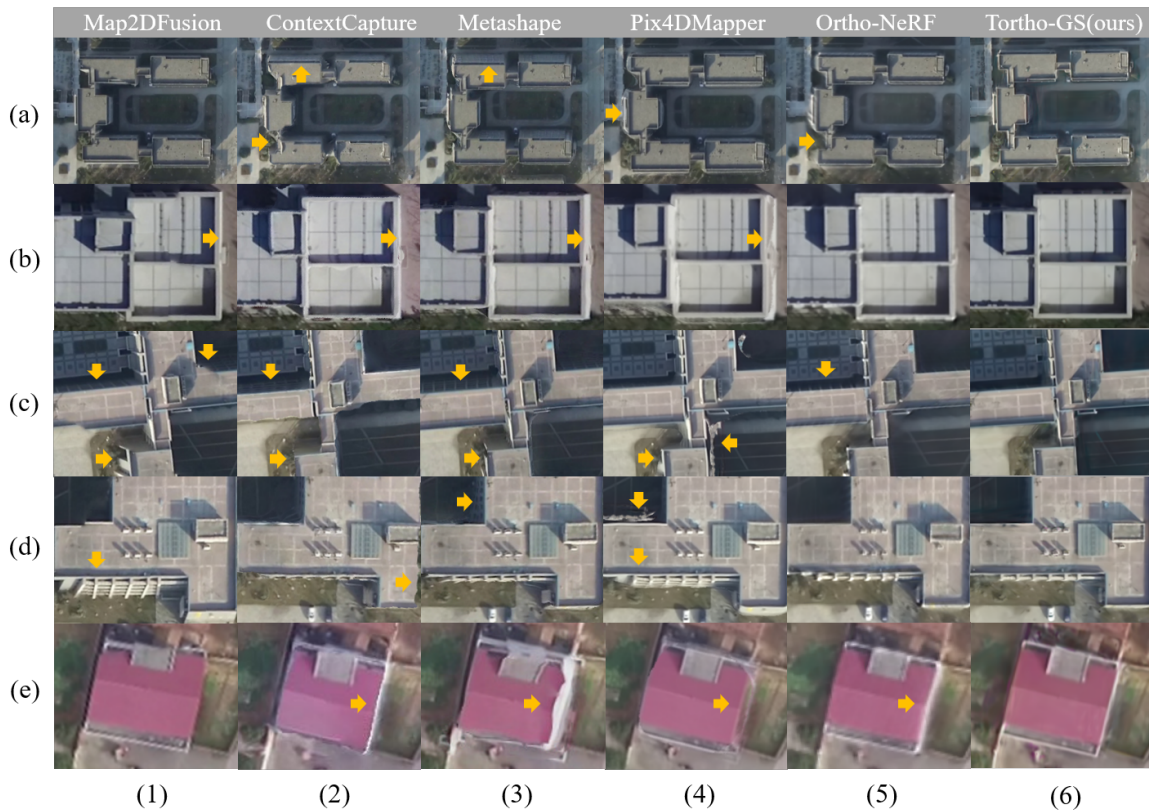
Fig. 6. Qualitative comparison of TDOMs generated by Map2DFusion, commercial software, Ortho-NeRF and our method on NPU DroneMap dataset. Our method can effectively reconstruct linear building edges and shows superior performance in eliminating distortions of building facades.

TABLE II
RELATIVE PRECISION OF TORTHO-GAUSSIAN COMPARING TO METASHAPE AND PIX4DMAPPER

| ID | Tortho-Gaussian | Metashape | | | Pix4DMapper | | |
|---|---|---|---|---|---|---|---|
| | Ratio | Ratio | Relative Error (%) | Absolute Error | Ratio | Relative Error (%) | Absolute Error |
| 1 | 1.41978 | 1.42071 | 0.06578 | 0.000935 | 1.42186 | 0.14623 | 0.001145 |
| 2 | 0.54305 | 0.54413 | 0.19844 | 0.001080 | 0.54253 | 0.09648 | 0.001603 |
| 3 | 1.01485 | 1.01614 | 0.12680 | 0.001288 | 1.01721 | 0.23185 | 0.001700 |
| 4 | 1.22322 | 1.22181 | 0.11548 | 0.001411 | 1.22285 | 0.02993 | 0.001405 |
| 5 | 0.81333 | 0.81238 | 0.11662 | 0.000947 | 0.81371 | 0.04705 | 0.001303 |
| 6 | 3.04875 | 3.04673 | 0.06639 | 0.002023 | 3.04333 | 0.17822 | 0.005424 |
| 7 | 1.04369 | 1.03974 | 0.37979 | 0.003499 | 1.04526 | 0.15076 | 0.001036 |
| 8 | 1.21988 | 1.22197 | 0.17061 | 0.001793 | 1.21833 | 0.12474 | 0.001557 |
| **Mean** | - | - | **0.15486** | **0.001622** | - | **0.12591** | **0.001772** |

prone to aliasing [72–74]. As the results of slender structures shown in Fig. 10 (a), (b), (c), (d), our method outperforms other approaches in both effectively and clearly restoring the geometric details of triangular power, towers, excavators and cranes. Our TOrtho-Gaussian approach demonstrates exceptional efficacy in reconstructing slender structures and thin shells, preserving the visual integrity of true orthographic images without distortion, breakage, or noise. This highlights its capability to handle challenging scenarios involving complex, fine-scale geometries with high fidelity.

*4) Weak Texture:* Weak texture often occurs in generating TDOM, such as water bodies and lake surfaces, which is dif-

ficult for traditional methods to fully reconstruct, as evidenced by the ghosting in Fig. 11 (a,4), (b,2) and Fig. 11 (c,1), the holes in Fig. 11 (a,2) and Fig. 11 (c,2), the tree reflections along the edges in Fig. 11 (b,1), (b,3), and (b,4). These are all results of failed weak-texture reconstructions, due to the limitations of conventional TDOM methods. Our Gaussian splatting technique offers an effective alternative to these challenges. By accurately fitting the Gaussian field to weak-texture regions, our method ensures a continuous, smooth surface, minimizing the case of holes. The differentiable nature of the Gaussian field further enhances this continuity, while true orthographic splatting effectively eliminates artifacts like

TABLE III
TRAINING TIME AND VRAM USAGE ON PHANTOM DATASET. BEST RESULTS ARE IN BOLD.

| Dataset | Phantom-ieu | | Phantom-factory | | Phantom-npu | |
|---|---|---|---|---|---|---|
| Method | Training | VRAM | Training | VRAM | Training | VRAM |
| Vanilla 3DGS | Failed | 21.6 GB | Failed | 22.8G | Failed | 22.2G |
| TOrtho-Gaussian | **39m23s** | **10.9G** | **42m5s** | **10.2G** | **43m7s** | **11.2G** |



Fig. 7. Qualitative comparison of TDOM generated by Map2DFusion, commercial software, and our method. The gray block denotes the result is not available due to the practical re-implementation issue.
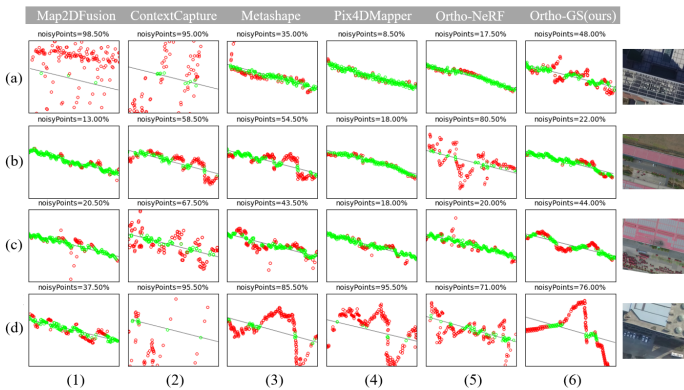


Fig. 8. Analysis of the building edges on various TDOM generated by Map2DFusion, commercial software, Ortho-NeRF and our method on NPU DroneMap dataset.
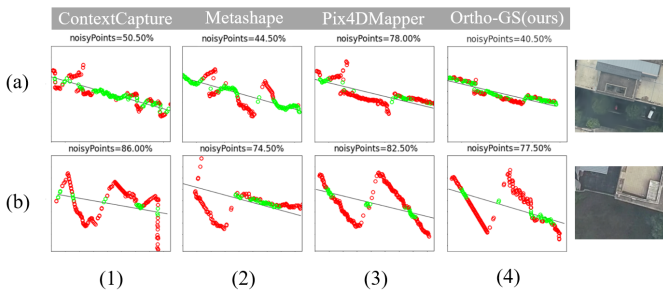


Fig. 9. Analysis of the building edges on various TDOM generated by commercial software, and our method on WHU dataset.
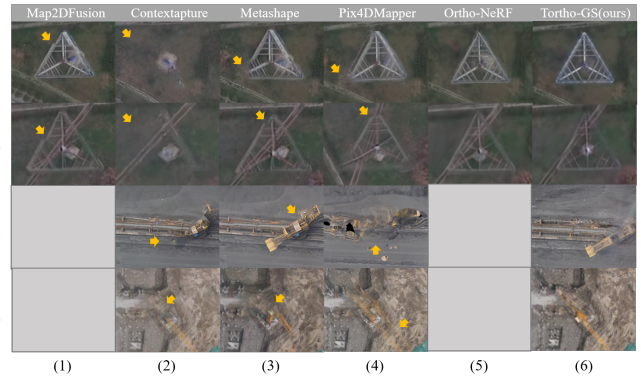


Fig. 10. Visual comparison of TDOMs using different methods on slender scenes. The gray block denotes that the result is not available due to the practical re-implementation issue
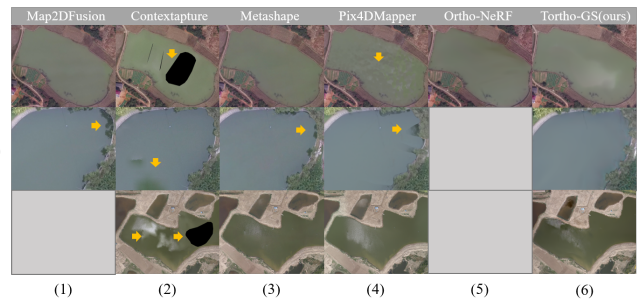


Fig. 11. Visual comparison of TDOMs using different methods on weak texture scenes. The gray block denotes that the result is not available due to the practical re-implementation issue

tree reflections. Additionally, the continuous Gaussian field acts as a filter, reducing erroneous texture noise and improving overall reconstruction quality.

### C. Quantitative Evaluation

Due to the absence of 3D control points, two popular commercial software, Metashape and Pix4DMapper, are applied as ground truth to evaluate the relative precision of our TDOM. We measure the lengths of line segments at building corner points on the TDOM from Tortho-Gaussian, Metashape, and Pix4DMapper. Each approach's TDOM has eight groups of two line segments, and the ratio of the two line lengths is taken as an indicator for TDOM quality, which ideally should remain identical across all DTOMs. Thus, we calculate the absolute and relative errors between these ratios as Tab. II lists. The average relative error and average absolute error of the ratios obtained from Tortho-Gaussian and Metashape are 0.155% and 0.001622, respectively. Similarly, the average relative

error and average absolute error between Tortho-Gaussian and Pix4DMapper are 0.126% and 0.001772, respectively. These findings indicate that the TDOM produced by our method achieves a level of accuracy comparable to state-of-the-art commercial software, demonstrating its reliability for precise geospatial reconstruction.

To further investigate the mapping accuracy, we conduct an overlay analysis using the Chengzi dataset by superimposing the TDOM generated by our method onto the corresponding CAD vector maps. As illustrated in Fig. 12, the overlay results demonstrate that our TDOM aligns precisely with the CAD maps, accurately replicating the shapes and boundaries of various features. This highlights the reliability and robustness of the Tortho-Gaussian approach in preserving geometric fidelity and ensuring accurate representation of spatial content.
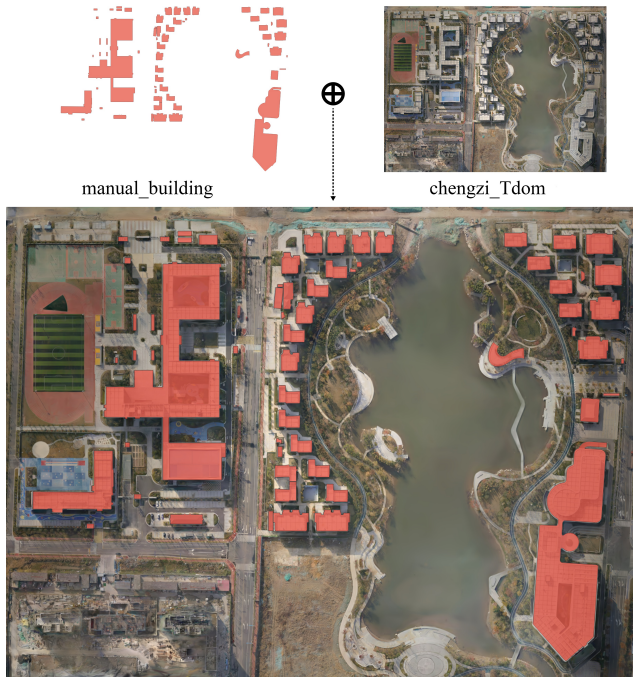


Fig. 12. The overlaid result with CAD map. Our generated TDOM is overlaid with an referenced CAD map produced via manual editing, the overlapping areas are highlighted with red background.

### D. Efficiency

We explored the optimization time and video memory cost for orthographic photos using the Tortho-Gaussian and vanilla 3DGS. Based on the Phantom3 dataset, both methods were tested with 30,000 iterations to ensure full convergence of the 3D Gaussian fields.

Tab.III provide the results from nine Gaussian field strips. Our proposed method can significantly reduces the cost time for reconstructing large-scale scenes. Notably, Tortho-Gaussian method consistently outperforms the vanilla 3DGS in terms of optimization time and memory efficiency when applied to large-scale scenes, demonstrating its practical value for commercial orthophoto production.

### E. Ablation Studies

In this section, we conduct in-situ ablation studies to look into various aspects of Tortho-Gaussian.

*1) Spatial Resolution:* As Section III-C explained that different spatial resolutions vary the resolution of TDOM and the number of 3D Gaussians for rendering a specific pixel. Thus, we test several set of $s_x$ and $s_y$ values to evaluate the influence of our orthogonal splatting at different spatial resolutions. Fig. 13 depicts that, in general, relevant TDOMs can be accurately rendered across different spatial resolutions, indicating that our Tortho-Gaussian is capable of generating TDOM products at multiple resolutions. Nevertheless, the fidelity and details of TDOM increase as the spatial resolution becomes higher and tend to be stable, this can be explained by the fact that, for large spatial resolution, more splat 2D Gaussians are input for $\alpha$-blending which may include some noisy and adjacent yet non-relevant Gaussians. As for the higher spatial resolution, more compact 2D Gaussians are considered and the fidelity of rendered TDOM tends to stable when the spatial resolution is approaching to the real spatial resolution of the input images.

Notably, the absence of FAGK in rendering lower spatial-resolution TDOM products leads to edge dilation artifacts, particularly along the boundaries of slender structures. This issue arises from the dilation and erosion effects of Gaussian ellipsoids during rendering, which introduce high-frequency Gaussian-shaped artifacts or pronounced swelling effects when the sampling rate (e.g., focal length or camera distance) is adjusted [71]. These artifacts can be mitigated through techniques such as multi-sampling, area sampling, or super-sampling. As illustrated in 13, the incorporation of FAGK effectively address this issue by introducing refined transparency transitions to the Gaussian kernels, thereby improving the sampling precision and significantly reducing edge-related artifacts.
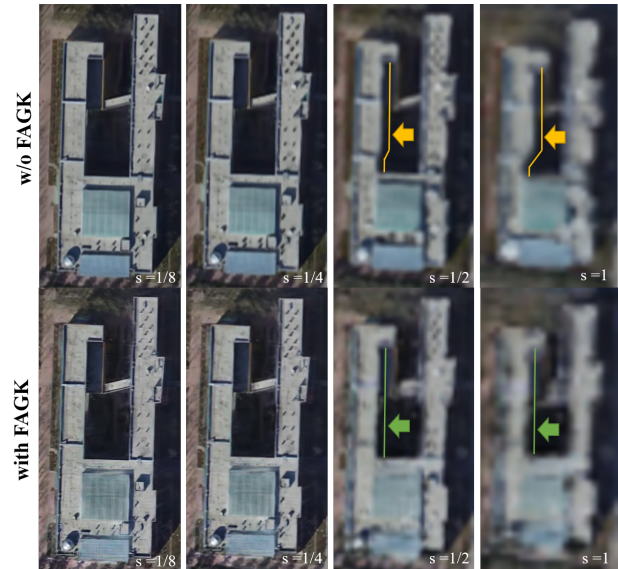


Fig. 13. Results of TDOM using different Spatial Resolutions(SR). From left to right, the spatial resolution decreases. As the screen space shrinks, the original 3DGS exhibits aliasing during projection onto 2D Gaussians, resulting in expansion artifacts that make the building structures appear thicker.

*2) Regional Partitioning:* Various partitions may result in different performance on the rendering results and cost time. Tab. IV compares different number of partitions, it can be found that increasing the number of partitioned regions generally results in a decline in SSIM, PSNR, and LPIPS values, indicating a degradation in image quality. This specific trend might due to that excessive partitioning can lead in additional boundary effects among various blocks. However, when the number of regions is set to one, there is an observable decrease in these values, which is most likely attributed to instable optimization and insufficient iterations for a whole large scene. As for the cost time, with four GPUs, 2×2 partition is the most time efficient, because this parallel optimization strategy leverages the full computational powers of multiple GPUs and each GPU can independently optimize a block in parallel, resulting in a substantial reduction in training time. The other partition solutions (with more than 4 partitions) need to wait free available computation resources (note that, as [37], the cost time can be further reduced if more multiple GPUs are employed to do parallel training for multiple blocks). These findings highlight the importance of balancing the level of partitioning to optimize image quality while minimizing the introduction of artifacts.

TABLE IV
ABLATION ON DATA PARTITION STRATEGY.

| Scene | Partition | Metrics | | | |
|---|---|---|---|---|---|
| | | SSIM | PSNR | LPIPS | Time |
| | 1 | 0.831 | 26.25 | 0.173 | 0h41m* |
| | 2×2 | **0.921** | **31.01** | 0.057 | **0h39m** |
| Phantom3-ieu | 3×3 | 0.919 | 30.77 | **0.053** | 1h35m |
| | 4×4 | 0.902 | 28.90 | 0.061 | 2h2m |
| | 5×5 | 0.907 | 28.95 | 0.059 | 3h19m |

* For this part, an modified 3DGS with optimized memory usage is applied. As section IV-A explains, 30000 iterations is run for the whole scene in total. And for 2×2 partitions, each subregion is also optimized by 30000 iterations.

*3) Fully Anisotropic Gaussian Kernel:* Tab. V provides the results of with/without the proposed FAGK using various degree setting (see section III-D). We can see that incorporating our fully anisotropic Gaussian kernels can typically improve the SSIM, PSNR, and LPIPS values, and higher degree can further enhance the rendering results. This improvement is attributed to the FAGK's ability to enable super-sampling, which enhances reflectivity sensitivity and provides better anti-aliasing effects.

TABLE V
ABLATION STUDY ON FULLY ANISOTROPIC GAUSSIAN KERNELS.

| Degree setting | Model setting | SSIM | PSNR | LPIPS |
|---|---|---|---|---|
| 1 | w/o FAGK | 0.813 | 18.22 | 0.285 |
| | Full, (Ours) | 0.806 | 19.27 | 0.296 |
| 2 | w/o FAGK | 0.867 | 23.18 | 0.223 |
| | Full, (Ours) | 0.867 | 23.80 | 0.224 |
| 3 (full) | w/o FAGK | 0.868 | 23.33 | **0.218** |
| | Full, (Ours) | **0.870** | **25.10** | 0.224 |

## V. CONCLUSION

In this paper, we presented a novel approach for generating True Digital Orthophoto Maps (TDOMs) by utilizing the 3D Gaussian Splatting (3DGS) technique, namely, Tortho-Gaussian. Our Tortho-Gaussian can bypass the traditional TDOM generation difficulties of occlusion detection and distortion correction, enabling a direct and efficient production of high-quality TDOM. Our throughout experimental results demonstrate that, Tortho-Gaussian not only outperforms existing commercial software in terms of TDOM precision and quality, but also provides a scalable solution for large-scale scene reconstruction. By dividing scenes into manageable segments and optimizing them independently, our method addresses the computational challenges associated with large datasets and achieves efficient memory usage without compromising rendering quality.

In the future, two potential relevant directions could be explored: first, to handle even more large-scale scenes (e.g., city-level), it is necessary to light the 3DGS and adopt a more reasonable partition strategy to minimize the thread wait among blocks. Second, we would like to incorporate other semantic information (e.g., segment anything model [75]) and depth (depth anything [76]) into our Tortho-Gaussain to further improve the quality of TDOM.

## REFERENCES

[1] S. S. Deshpande, "Bank line extraction by integration of orthoimages and lidar digital elevation model using principal component analysis and alpha matting," *Photogrammetric Engineering & Remote Sensing*, vol. 90, no. 10, pp. 631–638(8), 2024.

[2] B. A. DeWitt and P. R. Wolf, *Elements of Photogrammetry (with Applications in GIS)*. McGraw-Hill Higher Education, 2000.

[3] Y. Liu, X. Zheng, G. Ai, Y. Zhang, and Y. Zuo, "Generating a high-precision true digital orthophoto map based on uav images," *ISPRS International Journal of Geo-Information*, vol. 7, no. 9, p. 333, 2018.

[4] T. Li, C. Jiang, Z. Bian, M. Wang, and X. Niu, "A review of true orthophoto rectification algorithms," in *IOP Conference Series: Materials Science and Engineering*, vol. 780, no. 2. IOP Publishing, 2020, p. 022035.

[5] Y. Chen, C. Wang, W. Wang, X. Zhang, and N. Chen, "Building shadow detection based on improved quick shift algorithm in gf-2 images," *Photogrammetric Engineering & Remote Sensing*, vol. 90, no. 8, pp. 493–502(10), 2024.

[6] F. Amhar, J. Jansa, C. Ries *et al.*, "The generation of true orthophotos using a 3d building model in conjunction with a conventional dtm," *International Archives of Photogrammetry and Remote Sensing*, vol. 32, pp. 16–22, 1998.

[7] J. Rau, N.-Y. Chen, and L. Chen, "Hidden compensation and shadow enhancement for true orthophoto generation," in *Proceedings of Asian Conference on Remote Sensing 2000*, 2000, pp. 4–8.

[8] K. Uchida and T. Doihara, "Triangle-based visibility analysis and true orthoimage generation," *Int. Arch. Photogramm. Remote Sens. Spatial Inf. Sci.*, 2004.

[9] G. Zhou, W. Chen, J. A. Kelmelis, and D. Zhang, "A comprehensive study on urban true orthorectification," *IEEE Transactions on Geoscience and Remote sensing*, vol. 43, no. 9, pp. 2138–2147, 2005.

[10] A. F. Habib, K.-I. Bang, C. Kim, and S. Shin, "True ortho-photo generation from high resolution satellite imagery," *Innovations in 3D Geo Information Systems*, pp. 641–656, 2006.

[11] K. Bang, A. F. Habib, S. Shin, and K. Kim, "Comparative analysis of alternative methodologies for true ortho-photo generation from high resolution satellite imagery," *ASPRS ANNUAL*, vol. 2007, 2007.

[12] L.-C. Chen, T.-A. Teo, J.-Y. Wen, and J.-Y. Rau, "Occlusion-compensated true orthorectification for high-resolution satellite images," *The Photogrammetric Record*, vol. 22, no. 117, pp. 39–52, 2007.

[13] A. F. Habib, E.-M. Kim, and C.-J. Kim, "New methodologies for true orthophoto generation," *Photogrammetric Engineering & Remote Sensing*, vol. 73, no. 1, pp. 25–36, 2007.

[14] R. Antequera, P. Andrinal, R. González, S. Breit, J. Delgado, J. Pérez, M. Ureña, and S. Molina, "Development of an integrated system of true orthorectification. the altais lrto system," in *The International Archives of the Photogrammetry, Remote Sensing Spatial Information Sciences*, vol. 37, 2008, pp. 253–258.

[15] K.-I. Bang and C.-J. Kim, "A new true ortho-photo generation algorithm for high resolution satellite imagery," *Korean Journal of Remote Sensing*, vol. 26, no. 3, pp. 347–359, 2010.

[16] Y. P. Kuzmin, S. A. Korytnik, and O. Long, "Polygon-based true orthophoto generation," in *XXth ISPRS Congress Proceedings*, 2004, pp. 12–23.

[17] S. Wang, D. Li, Z. Guo, and J. Zheng, "Method for information processing of shadows and occlusion on orthophotos," *Journal of Geomatics (Chinese)*, vol. 29, no. 4, pp. 1–4, 2004.

[18] G. Zhou and Y. Wang, "Occlusion detection for urban aerial true orthoimage generation," in *2016 IEEE International Geoscience and Remote Sensing Symposium (IGARSS)*. IEEE, 2016, pp. 3009–3012.

[19] C. L. Zitnick and T. Kanade, "A cooperative algorithm for stereo matching and occlusion detection," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 22, no. 7, pp. 675–684, 2000.

[20] J. Zhang, B. Xu, M. Sun, and Y. Zhang, "True orthoimage generation based on occlusion detection with tin," *Geomatics and Information Science of Wuhan University*, vol. 37, no. 3, pp. 326–329, 2012.

[21] Y. H. Shin, S. W. Hyung, and D.-C. Lee, "True orthoimage generation from lidar intensity using deep learning," *Journal of the Korean Society of Surveying, Geodesy, Photogrammetry and Cartography*, vol. 38, no. 4, pp. 363–373, 2020.

[22] Y. H. Shin and D.-C. Lee, "True orthoimage generation using airborne lidar data with generative adversarial network-based deep learning model," *Journal of Sensors*, vol. 2021, no. 1, p. 4304548, 2021.

[23] S. A. N. Gilani, M. Awrangjeb, and G. Lu, "An automatic building extraction and regularisation technique using lidar point cloud data and orthoimage," *Remote Sensing*, vol. 8, no. 3, p. 258, 2016.

[24] M. Haggag, M. Zahran, and M. Salah, "Towards automated generation of true orthoimages for urban areas," *American Journal of Geographic Information System*, vol. 7, no. 2, pp. 67–74, 2018.

[25] Q. Wang, L. Yan, Y. Sun, X. Cui, H. Mortimer, and Y. Li, "True orthophoto generation using line segment matches," *The Photogrammetric Record*, vol. 33, no. 161, pp. 113–130, 2018.

[26] X. Li, H. Yan, S. Yang, and L. Niu, "A fusion algorithm of multispectral remote sensing image and aerial image," *Remote Sensing Information*, vol. 34, no. 4, pp. 11–15, 2019.

[27] G. Chen, S. Chen, X. Li, P. Zhou, and Z. Zhou, "Optimal seamline detection for orthoimage mosaicking based on dsm and improved jps algorithm," *Remote Sensing*, vol. 10, no. 6, p. 821, 2018.

[28] J. Yang, L. Liu, J. Xu, Y. Wang, and F. Deng, "Efficient global color correction for large-scale multiple-view images in three-dimensional reconstruction," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 173, pp. 209–220, 2021.

[29] G. Zhou, *Urban High-Resolution Remote Sensing: Algorithms and Modeling*. CRC Press, 2020.

[30] W. Yuan, X. Yuan, Y. Cai, and R. Shibasaki, "Fully automatic dom generation method based on optical flow field dense image matching," *Geo-spatial Information Science*, vol. 26, no. 2, pp. 242–256, 2023.

[31] Y. Liu, C. Luo, L. Fan, N. Wang, J. Peng, and Z. Zhang, "Citygaussian: Real-time high-quality large-scale scene rendering with gaussians," in *European Conference on Computer Vision*. Springer, 2025, pp. 265–282.

[32] S. Fridovich-Keil, A. Yu, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa, "Plenoxels: Radiance fields without neural networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 5501–5510.

[33] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan, "Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 5855–5864.

[34] L. Liu, J. Gu, K. Zaw Lin, T.-S. Chua, and C. Theobalt, "Neural sparse voxel fields," *Advances in Neural Information Processing Systems*, vol. 33, pp. 15 651–15 663, 2020.

[35] A. Yu, R. Li, M. Tancik, H. Li, R. Ng, and A. Kanazawa, "Plenoctrees for real-time rendering of neural radiance fields," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5752–5761.

[36] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman, "Mip-nerf 360: Unbounded anti-aliased

neural radiance fields," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 5470–5479.

[37] J. Lin, Z. Li, X. Tang, J. Liu, S. Liu, J. Liu, Y. Lu, X. Wu, S. Xu, Y. Yan *et al.*, "Vastgaussian: Vast 3d gaussians for large scene reconstruction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 5166–5175.

[38] Y. Sheng, "Minimising algorithm-induced artefacts in true ortho-image generation: a direct method implemented in the vector domain," *The Photogrammetric Record*, vol. 22, no. 118, pp. 151–163, 2007.

[39] E. T. Slonecker, B. Johnson, and J. McMahon, "Automated imagery orthorectification pilot," *Journal of Applied Remote Sensing*, vol. 3, no. 1, p. 033552, 2009.

[40] M. Ø. Nielsen, "True orthophoto generation," Master's thesis, Technical University of Denmark, DTU, DK-2800 Kgs. Lyngby, Denmark, 2004.

[41] Y. Sheng, P. Gong, and G. S. Biging, "True orthoimage production for forested areas from large-scale aerial photographs," *Photogrammetric Engineering & Remote Sensing*, vol. 69, no. 3, pp. 259–266, 2003.

[42] X. Wang, W. Jiang, and J. Xie, "A new method for true orthophoto generation," *Geomatics and Information Science of Wuhan University*, vol. 34, no. 10, pp. 1250–1254, 2009.

[43] H. Oliveira and M. Galo, "Occlusion detection by height gradient for true orthophoto generation, using lidar data," *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 40, pp. 275–280, 2013.

[44] C. Zhong, X. Huang, and D. Li, "Polygon-based inversion imaging for occlusion detection in true orthophoto generation," *Acta Geodaetica et Cartographica Sinica*, vol. 39, no. 1, p. 59, 2010.

[45] W. Xie, "Study on urban large-scale true orthophoto generation," *Acta Geodaetica et Cartographica Sinica*, vol. 39, no. 1, p. 0, 2010.

[46] C. Zhong, H. Li, Z. Li, and D. Li, "A vector-based backward projection method for robust detection of occlusions when generating true ortho photos," *GIScience & Remote Sensing*, vol. 47, no. 3, pp. 412–424, 2010.

[47] F. Deng, P. Li, Y. Kan, J. Kang, and F. Wan, "Overall projection of dbm for occlusion detection in true orthophoto generation," *Geomatics and Information Science of Wuhan University*, vol. 42, no. 1, pp. 97–102, 2017.

[48] Y. Hu, D. Stanley, and Y. Xin, "True ortho generation of urban area using high resolution aerial photos," *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 3, pp. 3–10, 2016.

[49] M. Ebrahimikia and A. Hosseininaveh, "True orthophoto generation based on unmanned aerial vehicle images using reconstructed edge points," *The Photogrammetric Record*, vol. 37, no. 178, pp. 161–184, 2022.

[50] M. Ebrahimikia, A. Hosseininaveh, and M. Modiri, "Orthophoto improvement using urban-snowflakenet," *Applied Geomatics*, vol. 16, no. 2, pp. 387–407, 2024.

[51] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[52] J. Lv, G. Jiang, W. Ding, and Z. Zhao, "Fast digital orthophoto generation: A comparative study of explicit and implicit methods," *Remote Sensing*, vol. 16, no. 5, p. 786, 2024.

[53] S. Chen, Q. Yan, Y. Qu, W. Gao, J. Yang, and F. Deng, "Ortho-nerf: generating a true digital orthophoto map using the neural radiance field from unmanned aerial vehicle images," *Geo-spatial Information Science*, pp. 1–20, 2024.

[54] G. Chen and W. Wang, "A survey on 3d gaussian splatting," *arXiv preprint arXiv:2401.03890*, 2024.

[55] A. Dalal, D. Hagen, K. G. Robbersmyr, and K. M. Knausgård, "Gaussian splatting: 3d reconstruction and novel view synthesis, a review," *IEEE Access*, 2024.

[56] B. Fei, J. Xu, R. Zhang, Q. Zhou, W. Yang, and Y. He, "3d gaussian as a new vision era: A survey," *arXiv preprint arXiv:2402.07181*, 2024.

[57] T. Wu, Y.-J. Yuan, L.-X. Zhang, J. Yang, Y.-P. Cao, L.-Q. Yan, and L. Gao, "Recent advances in 3d gaussian splatting," *Computational Visual Media*, vol. 10, no. 4, pp. 613–642, 2024.

[58] R. J. Cotton and C. Peyton, "Dynamic gaussian splatting from markerless motion capture reconstruct infants movements," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 60–68.

[59] G. Wu, T. Yi, J. Fang, L. Xie, X. Zhang, W. Wei, W. Liu, Q. Tian, and X. Wang, "4d gaussian splatting for real-time dynamic scene rendering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 20 310–20 320.

[60] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering." *ACM Trans. Graph.*, vol. 42, no. 4, pp. 139–1, 2023.

[61] M. Zwicker, H. Pfister, J. Van Baar, and M. Gross, "Ewa volume splatting," in *Proceedings Visualization, 2001. VIS'01.* IEEE, 2001, pp. 29–538.

[62] ——, "Ewa splatting," *IEEE Transactions on Visualization and Computer Graphics*, vol. 8, no. 3, pp. 223–238, 2002.

[63] T. Lu, M. Yu, L. Xu, Y. Xiangli, L. Wang, D. Lin, and B. Dai, "Scaffold-gs: Structured 3d gaussians for view-adaptive rendering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 20 654–20 664.

[64] Y. Chen and G. H. Lee, "Dogaussian: Distributed-oriented gaussian splatting for large-scale 3d reconstruction via gaussian consensus," *arXiv preprint arXiv:2405.13943*, 2024.

[65] H. Xie, Z. Chen, F. Hong, and Z. Liu, "Gaussiancity: Generative gaussian splatting for unbounded 3d city generation," *arXiv preprint arXiv:2406.06526*, 2024.

[66] S. Bu, Y. Zhao, G. Wan, and Z. Liu, "Map2dfusion: Real-

time incremental uav image mosaicing based on monocular slam," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 4564–4571.

[67] Bentley Systems, "ContextCapture Viewer," [Online]. Available: https://www.bentley.com/software/contextcapture-viewer/, Sep. 2022.

[68] Agisoft LLC, "Metashape," [Online]. Available: https://www.agisoft.com/, Sep. 2022.

[69] Pix4D, "Pix4DMapper," [Online]. Available: https://www.pix4d.com/product/pix4dmapper-photogrammetry-software/, Sep. 2022.

[70] J. Canny, "A computational approach to edge detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, no. 6, pp. 679–698, 1986.

[71] Z. Yu, A. Chen, B. Huang, T. Sattler, and A. Geiger, "Mip-splatting: Alias-free 3d gaussian splatting," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19 447–19 456.

[72] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman, "Zip-nerf: Anti-aliased grid-based neural radiance fields," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 19 697–19 705.

[73] W. Hu, Y. Wang, L. Ma, B. Yang, L. Gao, X. Liu, and Y. Ma, "Tri-miprf: Tri-mip representation for efficient anti-aliasing neural radiance fields," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 19 774–19 783.

[74] Y. Wang, J. Wang, Y. Qu, and Y. Qi, "Rip-nerf: learning rotation-invariant point-based neural radiance field for fine-grained editing and compositing," in *Proceedings of the 2023 ACM International Conference on Multimedia Retrieval*, 2023, pp. 125–134.

[75] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, "Segment anything," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 4015–4026.

[76] L. Yang, B. Kang, Z. Huang, X. Xu, J. Feng, and H. Zhao, "Depth anything: Unleashing the power of large-scale unlabeled data," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 10 371–10 381.

**Wendi Zhang** received his bachelor's degree in Geomatics Engineering from China University of Petroleum in 2023 and is currently pursuing a master's degree at Wuhan University, Wuhan, China. His research interests include the application of machine learning and computer vision in photogrammetry and implicit reconstruction.



**Associated Prof. Hong Xie** received the B.S., M.S., and Ph.D. degrees in photogrammetry and remote sensing from Wuhan University, Wuhan, China, in 2007, 2009, and 2013, respectively. He is currently an Associate Professor with the School of Geodesy and Geomatics, Wuhan University. His research interests include target detection based on image deep learning, point cloud data quality improvement, point cloud information extractionand model reconstruction, mobile mapping, and surveying.



**Dr. Haibin Ai** received his Ph.D. in Photogrammetry and Remote Sensing from Wuhan University in 2009. He is currently a Researcher at the Chinese Academy of Surveying and Mapping, serving as the Head of the Aerial and Space-borne Photogrammetry Research Group. He leads a team dedicated to research in digital photogrammetry for aerial and space borne platforms, remote sensing image processing, computer vision, and 3D reconstruction.



**Prof. Qiangqiang Yuan** (IEEE Member) received the B.S. degree in surveying and mapping engineering and the Ph.D. degree in photogrammetry and remote sensing from Wuhan University, Wuhan, China, in 2006 and 2012, respectively. In 2012, he joined the School of Geodesy and Geomatics, Wuhan University, where he is a Professor. He has published more than 100 research articles, including more than 80 peer-reviewed articles in international journals, such as Nature Communications, Remote Sensing of Environment, ISPRS Journal of Photogrammetry and Remote Sensing, IEEE TRANSACTIONS ON IMAGE PROCESSING, and IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING. His research interests include image reconstruction, remote sensing image processing and application, and data fusion. Dr. Yuan was a recipient of the Youth Talent Support Program of China in 2019, the Top-Ten Academic Star of Wuhan University in 2011, and the recognition of Best Reviewer of IEEE GEOSCIENCE AND REMOTE SENSING LETTERS in 2019. In 2014, he received the Hong Kong Scholar Award from the Society of Hong Kong Scholars and China National Postdoctoral Council. He is an associate editor of five international journals and has frequently served as a referee for more than 40 international journals for remote sensing and image processing.
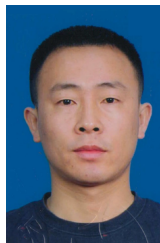


**Dr.-Ing Xin Wang** (IEEE Member) received his bachelor and master in surveying and mapping from school of Geodesy and Geomatics, Wuhan University, China, in 2013 and 2016, respectively. In 2021, he obtained Doctor of Engineering in photogrammetry and remote sensing from Leibniz university Hannover, Germany. He is currently an assistant professor in Wuhan University, whose research interests are computer vision, deep learning in applied photogrammetry etc.

**Prof. Zongqian Zhan** (IEEE Member) received the M.A.Eng. and Ph.D. degrees in photogrammetry and remote sensing from Wuhan University, Wuhan, China, in 2003 and 2007, respectively. He is currently a full Professor with the School of Geodesy and Geomatics, Wuhan University. His research interests include camera calibration, close-range and UAV photogrammetry, oblique photogrammetry, deep learning, and remote sensing.