

# Fast Satellite Tensorial Radiance Field for Multi-date Satellite Imagery of Large Size

Tongtong Zhang, Yuanxiang Li\*, Member, IEEE

arXiv:2309.11767v1 [cs.CV] 21 Sep 2023

**Abstract**—Existing NeRF models for satellite images suffer from slow speeds, mandatory solar information as input, and limitations in handling large satellite images. In response, we present SatensoRF, which significantly accelerates the entire process while employing fewer parameters for satellite imagery of large size. Besides, we observed that the prevalent assumption of Lambertian surfaces in neural radiance fields falls short for vegetative and aquatic elements. In contrast to the traditional hierarchical MLP-based scene representation, we have chosen a multiscale tensor decomposition approach for color, volume density, and auxiliary variables to model the lightfield with specular color. Additionally, to rectify inconsistencies in multi-date imagery, we incorporate total variation loss to restore the density tensor field and treat the problem as a denoising task. To validate our approach, we conducted assessments of SatensoRF using subsets from the spacenet multi-view dataset, which includes both multi-date and single-date multi-view RGB images. Our results clearly demonstrate that SatensoRF surpasses the state-of-the-art Sat-NeRF series in terms of novel view synthesis performance. Significantly, SatensoRF requires fewer parameters for training, resulting in faster training and inference speeds and reduced computational demands.

**Index Terms**—Neural Radiance Field, tensor decomposition, acceleration, DSM.

## I. INTRODUCTION

THE reconstruction of the earth's surface from high-resolution satellite imagery has consistently held significant importance across various domains. Conventional approaches within this domain typically tackle the challenge by solving multiple stereo problems to estimate disparity [1]–[5]. However, these approaches encounter inherent challenges during the dense matching phase. Despite the incorporation of deep learning techniques into satellite imagery, such as differentiable warping of Rational Polynomial Camera (RPC) for multi-view cost volume in [6], the computation and time costs are still substantial. Besides, most of the deep-learning based stereo problem require ground truth depth map as supervision [?].

In recent years, the emergence of the Neural Radiance Field (NeRF) [7] has sparked a revolutionary shift, particularly within the fields of computer graphics and computer vision. Reconstruction methods based on NeRF has much less dependency on dense matching thus requiring no groundtruth depth map for supervision. The fundamental concept behind the original NeRF is to use a simple, fully connected neural network to model a continuous scene.

Tongtong Zhang, Yuanxiang Li are with the Department of Information and Control, School of Aeronautics and Astronautics, Shanghai Jiao Tong University.

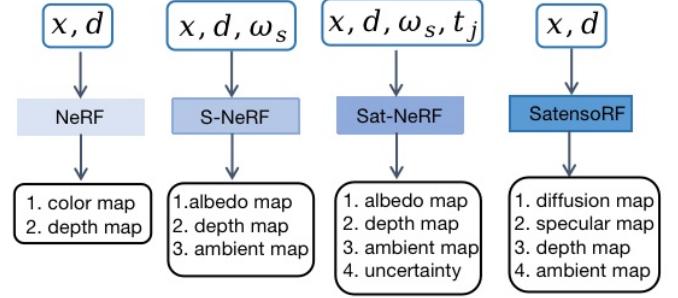


Fig. 1. The paradigm comparison of SatensoRF and the previous satellite-related NeRFs. SatensoRF has no requirements on additional input other than position  $x$  and viewing direction  $d$ , which is calculated from the sampled point cloud w.r.t. RPC.  $\omega_s$  refers to solar direction and  $t_j$  refers to the image index embeddings.

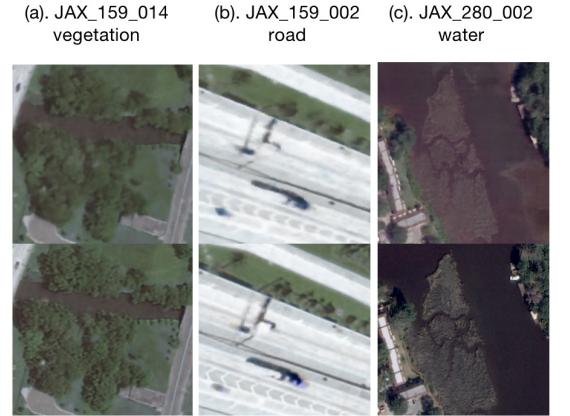


Fig. 2. Rendering results on specific patches with different light field assumptions. (a). For vegetational regions, SatensoRF with BRDF assumptions yields better tree details. (b). For reflective roads, SatensoRF results in finer edges. (c). SatensoRF is more capable of capturing water regions.

This network takes ray samples and viewing directions as inputs and maps them to RGB color and volume density. Subsequently, differential volume rendering facilitates end-to-end training for scene representation. Subsequent adaptations of NeRF in the aerospace domain have yielded significant benefits for drone applications [8], [9], as well as satellites [10], [11]. The common domain gap for migrating canonical RGB input to satellite images includes the following aspects: (1). radiometric inconsistencies among images across multiple dates. (2). discrepancies in the ray modeling between perspective cameras and rational polynomial cameras (RPCs), which require distinct sampling strategies. For the first gap, S-NeRF [10] introduces

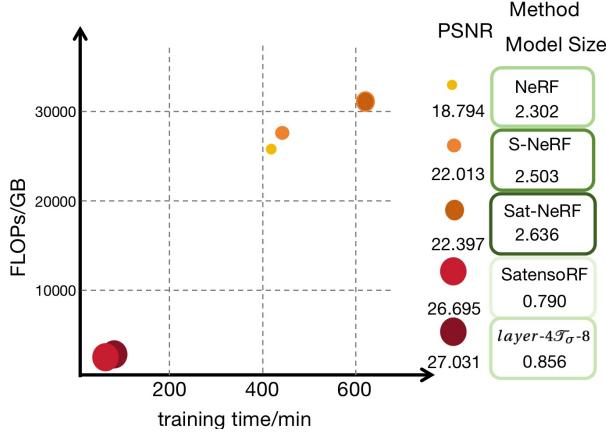


Fig. 3. Comprehensive comparison of SatensoRF with SOTA satellite NeRFs, on subset 159 of spacenet multi-view dataset with uncropped size. SatensoRF achieves better rendering results with faster training speed with less computation cost, and *layer-4T<sub>σ</sub>-8* refers to SatensoRF with 4 layers MLP and 8 levels of tensor radiance fields.

a remodeled neural volume rendering approach based on solar irradiance and incorporates the solar directions  $\omega_s$  as input, as depicted in (b) of Fig. 1. Building upon this, Sat-NeRF [11] extends the framework to model transient objects by regressing uncertainty to dynamically adjust the loss. Sat-NeRF also draws inspiration from NeRF in the wild (NeRF-w) [12], by introducing image index embedding  $t_j$  to supplement the multi-date difference in (c) of Fig. 1. For the second gap, Sat-NeRF also adapts the point sampling model to accommodate RPCs. However, despite the successful adaptation of neural radiance fields to multi-date satellite imagery, significant practical challenges persist.

The first left-over problem is that the conventional NeRF parameterization techniques in S-neRF and Sat-NeRF inherently retain the efficiency limitations of vanilla NeRF. In particular, the two-stage coarse-to-fine training approach requires a wide array of images from distinct viewpoints to ensure the seamless synthesis of continuous views. This rendering configuration results in dense querying of Multi-Layer Perceptrons (MLPs) at all query points. Besides, the purely MLP-based model introduces a trade-off between memory consumption and batch size, consequently impeding the overall efficiency and constraining the input image size.

The second problem is that Sat-NeRF requires additional input such as sun azimuth, elevation, as shown in Fig. 1. On the one hand, these additional input is essential for Sat-NeRF but might be unaccessible for some datasets. On the other hand, additional input tensors cause more computations, and further slower the running process. Specifically, we eliminate the following superfluous input variables:

- 1) The solar direction in [10], [11] used for forming the light field over the assumed Lambertian surface.
- 2) The index embedding in [10]–[12] to encode the transient objectives in multi-date photo collections.

The third problem is the over-simplified light field assumption. S-NeRF and Sat-NeRF have operated under the assumption that the earth's surface follows strict Lambertian law. They break down the rendered color into distinct elements—surface albedo, ambient color, and leveraging factor—dependent on solar orientation. Nevertheless, real-world surfaces often deviate from Lambertian qualities, exhibiting bidirectional reflection traits to varying degrees [13]. This phenomenon applies to surfaces like bodies of water, vegetation, and asphalt roads, and the distinction is shown in Fig.2.

For the first issue, our proposal tensorizes the 3D scene using Level-of-Detail representations. The appearance field of the earth's surface can be decomposed into multiscale representations, manifested as 3D voxel grids with three channels. Similarly, the volume density field of the Earth's surface and the global hue bias are expressed as 3D voxel grids with a single channel. SatensoRF largely improves the running speed and effectively reduces the total model size. Therein, SatensoRF requires less running memory and can be better applied to large satellite images than Sat-NeRF

For the second and third issue, we introduce a revised light field model with SatensoRF in Section III-C to account for the BRDF properties, without need of additional solar information and index embeddings. In Fig. 2, the proposed model encompassing the reflective light yields finer details of trees, metal surface of cars, and the water surface.

Moreover, regarding the cancelation of index embedding, we leverage the low-rank structure to handle the transient objects. This is achieved by treating transient objects as separate black-box noise components. Consequently, the task at hand transforms into the restoration of a stable volume density field, as further elucidated in Section III-D. Then, in virtue of the sparsity of the tensor radiance field, we apply total variation regularization to reduce the noise, and then prevents the influence of the inconsistency among images from different dates.

To sum up, our work makes several notable contributions:

- Introducing a lightweight and expedited tensor radiance field framework that requires less training and inference time.
- Reducing model sizes, and overall computational memory cost, and demonstrating superior results on large-sized satellite images, compared to previous NeRF-based approaches.
- Capturing the inherent anisotropic scattering properties of the earth's surface without presuming the Lambertian law, achieved through a lightweight model design with improved performance.
- Significantly reducing input dependencies by enabling SatensoRF to operate without relying on solar direction information.

The proposed SatensoRF are evaluated on different areas of interest the multi-date WorldView-3 images of the original

size. For altitude estimation, we adopt a lidar digital surface model (DSM) of resolution 0.5 m/pixel as evaluation reference. SatensoRF is compared to the state-of-the-art NeRF, S-NeRF, and Sat-NeRF.

## II. RELATED WORK

### A. Acceleration Framework for Neural Radiance Fields

In the conventional NeRF approach, the rendering step involves the dense evaluation of MLPs at all query points along each camera ray during numerical integration, leading to slow speed and extensive computational demands. To address this, baking techniques have been introduced to accelerate inference time [14]–[16]. In pursuit of swifter training and inference, the utilization of diverse formats for learnable features yields more concise outcomes. While MLPs exhibit competitiveness in neural implicit representations [17] and original NeRFs, they lag behind methods that directly optimize and query the scene representations themselves in terms of speed. Notably, Instant-NGP [18] employs multiresolution hash tables for scene representation organization, achieving fast processing with a level of detail. Other structures like trees [19], and sparse grids [15], [20] are also popular. In parallel, methods such as [21]–[23] structure the scene as tensor compositions with low rank. TensoRF [21], for instance, employs vector-matrix (VM) and CANDECOMP/PARAFAC (CP) decomposition to characterize the scene, followed by input to compact MLPs for efficient computation. In a similar vein, CCNeRF [22] discards the use of MLPs to facilitate compression and composition. NRFF extends TensoRF into a multiscale framework to enable comprehensive learning at varying scales, thereby enhancing accuracy.

### B. Neural Radiance Field with Multi-Date Images

The vanilla NeRF model conditions on the fact that the 3D scene have the density, radiance and illumination, which is unreal for outdoor scenes, especially for internet photo collections and satellite images. Both of these two types of dataset have different lighting condition, appearance of the pixel at the same 3D position. Besides, there are transient occlusions from people and cars. NeRF-W [12] firstly proposes a per-image latent embedding capturing photometric appearance variations often present in wild heritage dataset, and explicitly model the scene into image-dependent and shared components respectively, to disentangle transient elements from the static scene. With the same starting point, Sat-NeRF adapt the method to satellite images by representing the albedo of the earth surface, and the other lighting factors in different heads of MLPs. NeRF-OSR [24] enables outdoor scene relighting

### C. Deep Learning in Reconstruction from Satellite Images

Traditional 3D reconstruction has predominantly followed a multi-view stereo (MVS) pipeline that generates a dense disparity map using dense matching algorithms such as Semi-Global Matching (SGM) [25]. However, errors tend

to accumulate from the initial stereo pair selection to the subsequent refinement steps [25], [26]. As deep learning methodologies infiltrate the realm of 3D reconstruction, an increasing number of end-to-end differentiable MVS pipelines have emerged [27]. Some approaches extract features from the 3D cost volume and utilize 3D Convolutional Neural Networks (CNNs) [28], [29], Recurrent Neural Networks (RNNs) [30], or attention modules [31] to predict the disparity map. However, these methods often require significant computation and time. On the other hand, some methods adapt randomized iterative patch-match techniques [32], [33] into an end-to-end framework [34], [35].

The proliferation of Implicit Neural Representations (INR) introduces increased flexibility in terms of computation and memory utilization. Traditional epipolar geometry and INR complement each other in the realm of 3D reconstruction. MVS provides geometric guidance for the radiance field [36]–[38], while INR reduces the reliance on conventional MVS pipelines for selecting optimal stereo pairs and explicit depth supervision [39], [40].

## III. METHOD

### A. Light Field over the Earth Surface

The primary objective of SatensoRF is to generate a synthesized novel view photo in RGB, along with an altitude map product. The input of SatensoRF include the spatial coordinates of the scene volume, and the viewing direction  $\mathbf{d}$ . For a point  $\mathbf{x}_i = \mathbf{r}(t_i) = \mathbf{o} + t_i\mathbf{d}$ , the proposed tensor radiance field yields both the rendered color  $\mathbf{c}(\mathbf{r})$  and the rendered altitude  $\mathbf{h}(\mathbf{r})$  for the ray  $\mathbf{r}$  at a given sample  $t_i$  within the interval  $[t_{min}, t_{max}]$ . Here,  $t_{min}$  and  $t_{max}$  define the sampling range along the altitude direction.

The radiance observed by satellites over the atmosphere is the backscattered portion of incident solar radiation, which is composed of the pure backscattering and the radiance reflected by the underlying surface and transmittance to the satellite sensors. The three principal components are illustrated in Fig. 4:

- 1) **ambient radiation**  $L_{ambient}$ : This pertains to photons that are scattered into the satellite's optical sensor without encountering the ground surface directly (depicted in gray).
- 2) **specular radiation**  $L_{specular}$ : This denotes the specular components of the reflected radiation.
- 3) **diffusion radiation**  $L_{diffuse}$ : This involves the diffuse component of the reflected light.

In this context, it's essential to emphasize that the ambient light component  $L_{ambient}$  remains independent of the characteristics of the ground surface. Furthermore, the intricacies involved in tracing rays for scattered light  $L_{diffuse}$  prompts us to omit the calculation of parameters for  $L_{ambient}$  and  $L_{diffuse}$  across every volume within the scene. Instead, we simplify these aspects through a unified irradiance model. Consequently, these factors blend in

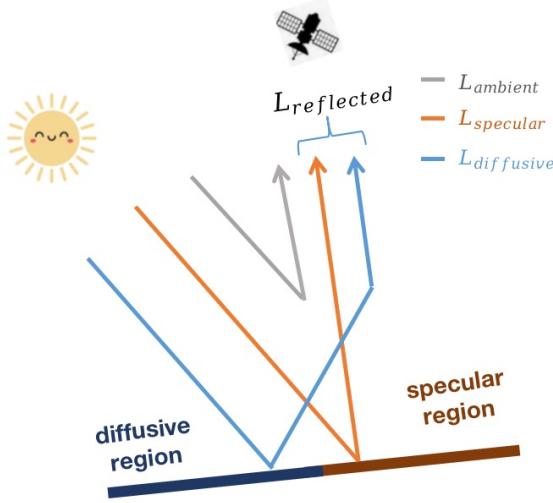


Fig. 4. Sketch of the radiance components for the ground surface.

the overall radiation signal, which can be summarized as follows:

$$\begin{aligned} L_{all} &= L_{ambient} + L_{diffuse} + L_{specular} \\ &\triangleq L_{ambient} + L_{reflected} \end{aligned}$$

In contrast to S-NeRF and Sat-NeRF with heavy reliance on the solar direction, our approach resort to a dependent randomly initialized low-rank tensor field to decouple from the earth surface color. Instead, we rely on another small tensor field serving as the latent code randomly sampled from normal distribution, and train it to represent the ambient map. Ultimately, for the 3D point  $\mathbf{x}_i$ , we achieve the following composition:

$$\begin{aligned} \mathbf{l}(\mathbf{x}_i) &= \lambda_{amb}(\mathbf{x}_i)\mathbf{l}_3 + (1 - \lambda_{amb}(\mathbf{x}_i))\mathbf{c}_{amb}(\mathbf{x}_i) \quad (1) \\ \mathbf{c}(\mathbf{x}_i, \mathbf{d}) &= \mathbf{c}_{ref}(\mathbf{x}_i, \mathbf{d}) \odot \mathbf{l}(\mathbf{x}_i) \end{aligned}$$

where  $\lambda_{amb}$  is the factor to adjust the impact of ambient radiation  $L_{ambient}$  and the ambient radiation  $L_{diffuse}$ . The ambient color  $\mathbf{c}_{amb}(\mathbf{x}_i)$  at  $\mathbf{x}_i$  to form the irradiance map  $\mathbf{l}(\mathbf{x}_i)$  to model  $L_{diffuse}$ , the reflected field  $\mathbf{c}_{ref}$  to model  $L_{reflected}$  will be discussed in detail in Section III-B. So the final rendered map  $\hat{\mathbf{l}}$  is:

$$\begin{aligned} \hat{\mathbf{l}} &= \sum_{i=1}^{N_s} Tr(\mathbf{x}_i)\alpha(\mathbf{x}_i)\mathbf{c}(\mathbf{x}_i, \mathbf{d}) \quad (2) \\ &= \sum_{i=1}^{N_s} Tr(\mathbf{x}_i)\alpha(\mathbf{x}_i)\mathbf{c}_{ref}(\mathbf{x}_i, \mathbf{d})\mathbf{l}(\mathbf{x}_i) \end{aligned}$$

where the opacity  $\alpha(\mathbf{x}_i)$  and transmittance  $Tr(\mathbf{x}_i)$  are calculated as

$$\alpha(\mathbf{x}_i) = 1 - \exp(-\sigma(\mathbf{x}_i)\delta_i) \quad (3)$$

$$Tr(\mathbf{x}_i) = \prod_{j=1}^{i-1} (1 - \alpha(\mathbf{x}_j)) \quad (4)$$

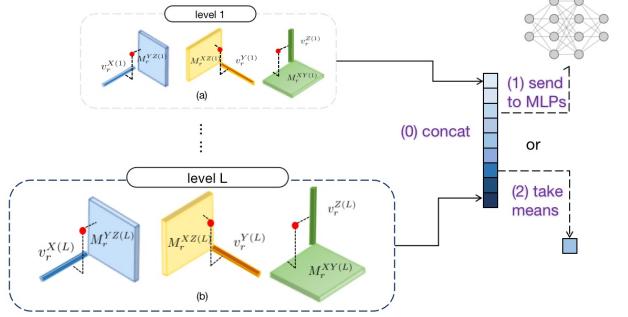


Fig. 5. Multiscale tensor factorization of the 3D scene before sending to the MLP

with  $\delta_i$  being the two consecutive points along the ray, i.e.  $\delta_i = t_{i+1} - t_i$  and  $N_s$  is the number of sample in the altitude direction.

### B. Scene as Tensor Decomposition in LoD

Within our tensor radiance field, the atomic representation involves a singular level of VM decomposition, which factorizes a tensor into multiple vectors and matrices. This structure is illustrated in (a) of Fig. 5. We designate the three dimensions of the 3D space as  $XYZ$ , wherein  $X$  and  $Y$  correspond to the north and east directions in the Universal Transverse Mercator (UTM) coordinate system, while  $Z$  signifies the elevation of the Earth's surface. Concretely, VM decomposition for a given single level  $l$  can be formulated as follows:

$$\mathcal{T}^{(l)} = \sum_{r=1}^R (v_r^{X(l)} \odot M_r^{YZ(l)} + v_r^{Y(l)} \odot M_r^{XZ(l)} + v_r^{Z(l)} \odot M_r^{XY(l)}) \quad (5)$$

where  $M_r^*$  are matrix factors for two of the three directions,  $\odot$  is the element-wise multiplication and  $R$  is the rank at three direction.

The coarse-to-fine strategy is a prevalent approach in NeRF pipelines. In our context, we have decomposed the feature field  $r$  for the 3D scenes into  $L$  levels. The sizes of each level are selected using a strategy akin to that of NGP [18], based on logarithmic considerations. Illustrations of the tensor factorization for level 1 and level  $L$  can be found in (a) and (b) of Fig. 5, respectively. It's important to note that distinct feature fields undergo varied procedures in subsequent stages.

To model both the radiance field and the irradiance field, SatensoRF employs feature fields that consist of several subfields, each with the following notations:

- The volume density field  $\mathcal{T}_\sigma$ : This represents the radiance feature field for volume density. It has  $L = 8$  levels and 4 feature channels. The outputs are concatenated and averaged across levels to form  $\sigma$ .
- The reflective field  $\mathcal{T}_{ref}$ : This represents the radiance feature field for appearance. It also has  $L = 8$  levels and

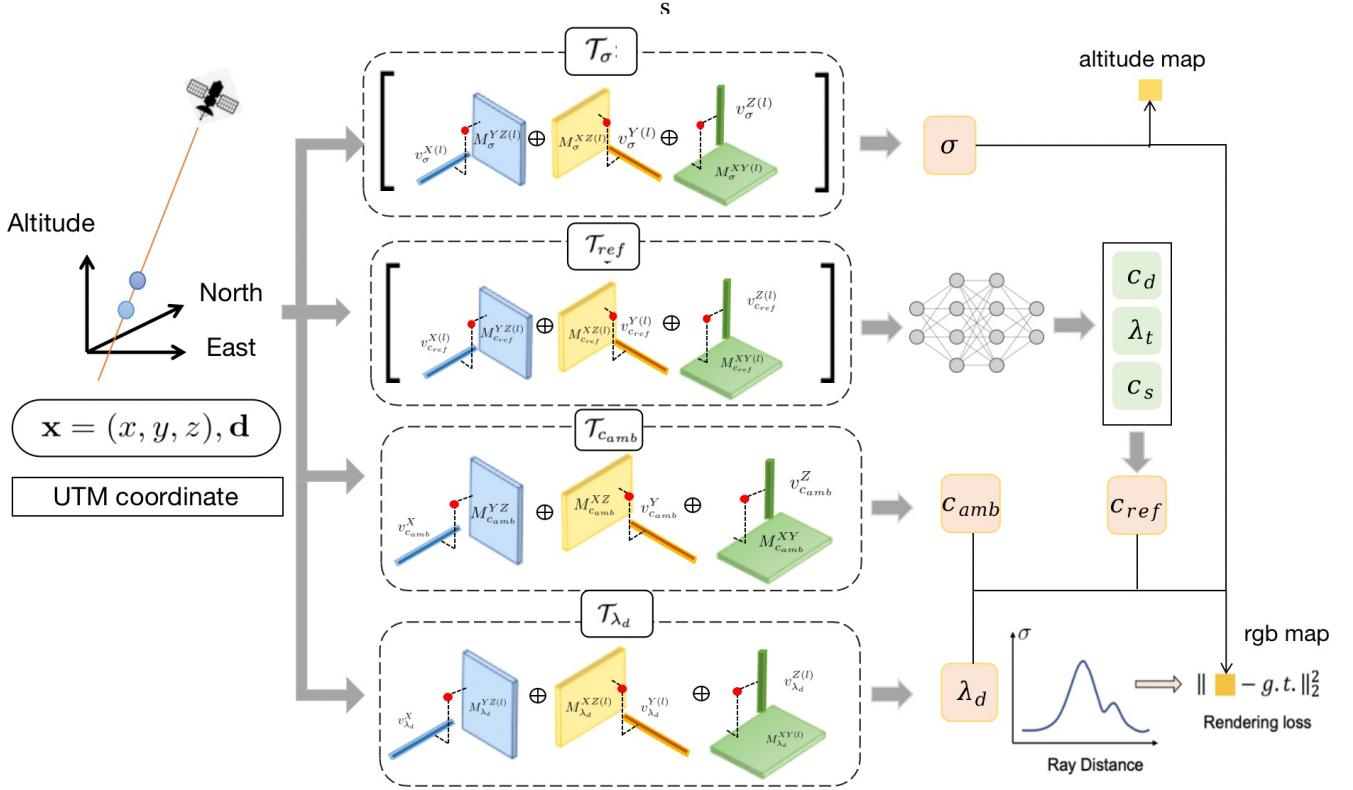


Fig. 6. The SatensoRF panorama encompasses four distinct sub-tensor radiance fields. These are outlined as follows: 1). The multi-scale volume density field  $\mathcal{T}_\sigma$  is computed and then averaged to produce the volume density field. 2). The multiscale reflectance field  $\mathcal{T}_{ref}$  is channeled into the MLPs for further processing. 3). Both the ambient color field  $\mathcal{T}_{camb}$  and the ambient ambient adjusting factor  $\mathcal{T}_{\lambda_{amb}}$  operate on a single-level scale. This is due to the fact that the finer details of the ground surface have minimal influence on ambient colors. After volume rendering, the loss function is applied to the synthesized image to facilitate the optimization process. This panorama demonstrates the various components and their interactions within the SatensoRF framework.

4 feature channels. The outputs are concatenated and fed into a small multi-head MLP network, as shown in Fig. 7. Further details will be provided in Section III-C.

- The ambient color field  $\mathcal{T}_{camb}$ : This represents the irradiance feature field for ambient color. It has only one level of tensor decomposition and consists of 3 feature channels, without multiscale details. The outputs are concatenated and averaged across levels to form  $\mathbf{c}_{amb}$ .
- The ambient factor field  $\mathcal{T}_{\lambda_{amb}}$ : This represents the irradiance field feature field. It has only one level of tensor decomposition and includes 2 feature channels, without multiscale details. The outputs are concatenated and averaged across levels to form  $\lambda_{amb}$ .

### C. The Reflective Radiance Field

The result color  $\mathbf{c}_{ref}$  of the reflective light  $L_{reflected}$  can be decomposed into diffusion color  $\mathbf{c}_d$  and specular light  $\mathbf{c}_s$  with factor  $\lambda_s$  [41], [42]:

$$\begin{aligned} \mathbf{c}_{ref}(\omega_o; \mathbf{x}) &= \mathbf{c}_d + \lambda_s \mathbf{c}_s \\ &= \mathbf{c}_d + \lambda_s \int_{\omega} L_i(\omega_i; \mathbf{x}) \rho_s(\omega_i, \omega_o; \mathbf{x}) (\mathbf{n} \cdot \omega_i) d\omega_i \\ &= \mathbf{c}_d + \lambda_s \int_{\omega} f(\omega_i, \omega_o; \mathbf{x}, \mathbf{n}) d\omega_i \end{aligned} \quad (6)$$

where  $L_i$  indicates the incoming radiance from direction  $\omega_i$  around the surface normal  $\mathbf{n}$  at position  $\mathbf{x}$ ,  $\rho_s$  is the bidirec-

tional reflectance distribution function (BRDF) depending on  $\mathbf{x}$  and  $\omega_i$ , and the dot product between  $\omega_i$  and  $\mathbf{n}$  serves as the cosine of the intersecting angle. The specular part is the integration over the hemisphere  $\Omega$  centered around  $\mathbf{n}$  containing all possible values for  $\omega_i$  where  $\omega_i \cdot \mathbf{n} > 0$ . Instead of solving the rendering equation via Monte Carlo estimation, implicit representations learn the components of the rendering equation via neural networks [23], [43].

Subsequently, we perform a decomposition of the reflective radiance field using the ambient color feature field  $\mathcal{T}_{camb}$ . We then feed the extracted features into Multi-Layer Perceptrons (MLPs). According to Equation 7, both the diffusion color  $c_d$  and the factor  $\lambda_s$  remain independent of the viewing direction. As a result, they can be directly regressed through MLPs. In this case, a 4-layer MLP structure (labeled as "A" in Fig. 7) is employed to predict the diffusion color  $c_d$ , the factor  $\lambda_s$ , and the normals  $n$  for every volume intersected by the sampled ray.

To improve the parameterization of the rendering equation, the encoding process is applied within the color feature space, which differs from most NeRF methods that perform encoding in the scene coordinate space [7], [23]. This approach is particularly advantageous because the BRDF function  $f$  depends solely on the specific features of the color domain. The encoded feature of the rendering

equation is modeled as an Anisotropic Spherical Gaussian (ASG) function, a concept introduced by the work [44]. This ASG function is composed of subfunctions  $G_i$  that further enhance the ability to represent intricate scene characteristics:

$$\begin{aligned}\mathcal{F}_s(\omega_o, \mathbf{x}) &= \sum_{i=0}^{N-1} G_i(\omega_o; \mathbf{x}, [\omega_i, \omega_i^\lambda, \omega_i^\mu], [\lambda_i, \mu_i], a_i) \\ &= \sum_{i=0}^{N-1} F_i \max(\omega_o \cdot \omega_i, 0) \exp(-\lambda_i (\omega_o \cdot \omega_i^\lambda)^2 - \mu_i (\omega_o \cdot \omega_i^\mu)^2)\end{aligned}\quad (7)$$

where  $F_i$  is the feature vector,  $\mathcal{B} = [\omega_i, \omega_i^\lambda, \omega_i, \omega_i^\mu]$  are orthonormal basis with  $\lambda_i, \mu_i > 0$  being the bandwidth at each basis.

The parameters of the ASG  $F_i$ ,  $\lambda_i$ , and  $\mu_i$ , are initially obtained from two MLPs. The first MLP labeled "A" (blue) and another 2-layer MLP labeled "B" (pink) in Fig. 7 contribute to determining these ASG parameters. The basis  $\mathcal{B}$  required for ASG representations is sampled on a unit sphere. These basis vectors are then processed through the encoder labeled "C" (green) in Fig. 7. The encoder generates the rendering equation that is contingent on the viewing direction, together with the intersection angle derived from the dot product of the normals  $n$  and the viewing direction  $d$ . Subsequently, the specular part  $\mathbf{c}_s$  is decoded using the last 2-layer MLP labeled "D" in Fig. 7, employing a Sigmoid activation function. This decoded specular part  $\mathbf{c}_s$  then contributes to forming the reflective color  $\mathbf{c}_{ref}$ .

#### D. Handling Transient Objects in Density Reconstruction

Instead of modeling the effects caused by transient objects by regressing the uncertainty, as proposed in [11] [11], we capitalize on the inherent advantages of low-rank tensor representations.

Total Variation Denoising (TVD) is a widely employed method for signal restoration tasks. This regularization technique is particularly effective in preserving sharp edges present within the underlying signal.

In this context, we treat the Earth's ground surface as the signal to be restored from the noisy input data. Meanwhile, transient objects such as cars or people are considered as the noise to be filtered out.

We denote the true volume density tensor as  $\mathcal{T}_\sigma$  and subsequently formulate the volume density tensor as follows:

$$\mathcal{T}_\sigma = \mathcal{X}_\sigma + \mathcal{N} \quad (8)$$

where  $\mathcal{N}$  is the noise tensor caused by the transient object with no prior rules of occurrence.

Remote sensing images often exhibit a prominent piecewise smooth structure, making them well-suited for the application of Total Variation (TV) methods in image processing [45], [46]. The underlying rationale of Total Variation regularization is rooted in the notion that excessive and spurious details tend to exhibit higher total variation [47], [48].

In light of this, we apply the concept of Total Variation loss to the decomposed planes of the tensor field  $\mathcal{T}_\sigma$  across all  $L$  levels. This approach is geared towards promoting smoother and more coherent structures within the data.

$$\mathcal{L}_{TV}(\mathcal{T}_\sigma) = \sum_{l=1}^L \|\mathbf{D}\mathcal{T}_\sigma\|_2 / L \quad (9)$$

$$\begin{aligned}\|\mathbf{D}\mathcal{T}_\sigma^{(l)}\|_2 &= \sum_{i=1}^L (\sum_j (M_\sigma^{XY(l)}(i+1,:) - M_\sigma^{XY(l)}(i,:)) \\ &\quad + \sum_j (M_\sigma^{XY(l)}(:,j+1) - M_\sigma^{XY(l)}(:,j)))\end{aligned}\quad (10)$$

where  $\mathbf{D}$  is the difference operator, and  $i, j$  indexes elements on the factorized matrix  $M_\sigma^{XY(l)}$  at scale  $l$  from the  $L$  levels.

#### E. The Losses

a)  $\mathcal{L}_{rgb}$ : The basic supervisions is the mean squared error (MSE) between the rendered color at every ray  $\mathbf{R}$  and the ground truth color of the input images:

$$\mathcal{L}_{rgb} = \sum_{\mathbf{r} \in \mathcal{R}} \|\hat{\mathbf{c}}(\mathbf{r}) - c_{gt}((\mathbf{r})\|_2^2 \quad (11)$$

b)  $\mathcal{L}_{normal}$ : Additionally we should force the estimated normals  $\mathbf{n}_i$  to be on the unit sphere:

$$\mathcal{L}_{normal} = \frac{1}{N} \sum_{i=0}^{N-1} \omega_i \max(0, \mathbf{d} \cdot \mathbf{n}_i)^2 \quad (12)$$

where  $N$  is the number of samples in the cast ray and  $\mathbf{d}$  is the viewing direction.

c)  $\mathcal{L}_{\lambda_{amb}}$ : The ambient light  $L_{ambient}$  takes no effect when the ray intersects with the earth surface, thus we encourage the ratio of the ambient light  $\lambda_{amb}$  to be close to the transmittance  $Tr$ .

$$\mathcal{L}_{\lambda_{amb}} = \sum_{i=1}^N (Tr(x_i) - \lambda_{amb}(x_i))^2 + 1 - \sum_{i=1}^N Tr(x_i) \alpha(x_i) \lambda_{amb}(x_i) \quad (13)$$

where  $N$  is number of total voxels.

d)  $\mathcal{L}_{depth}$ : To further enhance model performance, particularly concerning altitude reconstruction, SatensoRF incorporates depth supervision, akin to Sat-NeRF and DS-NeRF [38]. Utilizing precomputed sparse 3D points derived from SFM pipelines or bundle-adjusting procedures, the depth supervision is formulated as follows:

$$\mathcal{L}_{DS} = \sum_{\mathbf{r} \in \mathcal{R}_{DS}} w(\mathbf{r}) (\hat{H}(\mathbf{r}) - \|\mathbf{X}(\mathbf{r}) - \mathbf{o}(\mathbf{r})\|_2)^2 \quad (14)$$

where  $\mathcal{R}_{DS}$  is the set of rays containing the known 3D points,  $\|\mathbf{X}(\mathbf{r}) - \mathbf{o}(\mathbf{r})\|_2$  is the relative altitude as ground truth at ray  $\mathbf{r}$ ,  $\hat{H}(\mathbf{r})$  is the rendered height as estimation, and  $w(\mathbf{r})$  comes from the reprojection error from the bundle adjustment procedure <sup>1</sup>.

e) *Total loss*: The basic whole loss include the data fidelity term  $\mathcal{L}_{rgb}$ , the noise sparse prior  $\mathcal{L}_{TV}$ , the normal regularization  $\mathcal{L}_{normal}$  and the ambient light regularization  $\mathcal{L}_{\lambda_{amb}}$ :

$$\begin{aligned}\mathcal{L} = & \mathcal{L}_{rgb}(\mathbf{c}, c_{gt}) + \lambda_{TV} \mathcal{L}_{TV}(\mathcal{T}_\sigma) \\ & + \lambda_n \mathcal{L}_{normal} + \lambda_{\lambda_{amb}} \mathcal{L}_{\lambda_{amb}}\end{aligned}$$

<sup>1</sup><https://github.com/centreborelli/sat-bundleadjust>

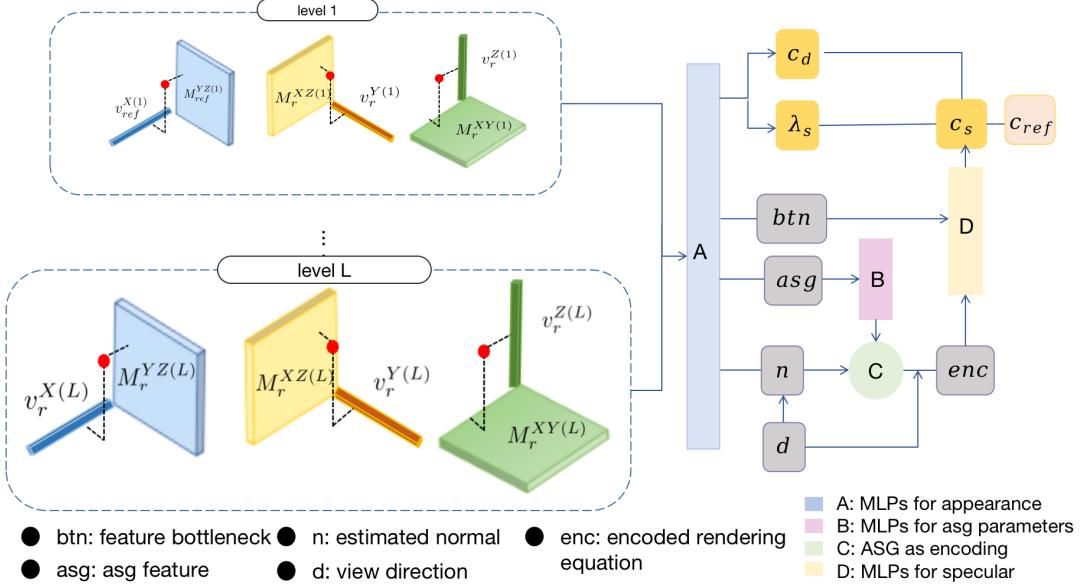


Fig. 7. The multihead MLPs to decouple the reflective light field.

#### IV. EXPERIMENTS

We evaluate the performance of SatensoRF based on the subsets of spacenet MVS-dataset<sup>2</sup>.

The evaluation is conducted from the following aspects.

- 1) For both single-date and multi-date images, we validate the effectiveness of SatensoRF, and compare it to the state-of-the-art Sat-NeRF and other stereo pipelines in Section IV-C and Section IV-B.
- 2) We validate the non-Lambertian assumption for SatensoRF in Section IV-D.
- 3) We conduct ablation study on total variation loss for SatensoRF in Section IV-E.
- 4) We conduct controled experiments to explore the effects of changing some variables in Section IV-F.

##### A. Setup

1) *Dataset Preparation:* Firstly for multi-date images, we try to validate that SatensoRF applies well to the large images with minimal available primitive information. So we choose the original dataset from the Maxar WorldView-3 satellite, capturing imagery between 2014 and 2016 over the city of Jacksonville, Florida, in the United States. In contrast to the multi-date images used by S-NeRF and Sat-NeRF, our experiments are conducted using full-size RGB images, each containing over 4 million pixels. This approach differs from cropping images into smaller patches based on areas of interest (AOI).

The original dataset<sup>3</sup> comprises multi-date imagery captured under various lighting conditions, and specific details

can be found in Table I. The original multi-date dataset contains approximately ten times as many rays as the DFC dataset, and it includes challenging scenarios. For instance, Subset 031 features a mix of vegetation, buildings, and asphalt roads, Subset 159 includes a reflective flyover, and Subsets 260 and 280 contain bodies of water.

Notably, Not all of the subsets of Jacksonville has corresponding multispectral images (MSI) to offer the extra information about sun elevation and sun azimuth, so we pick up two subsets with extra sun information (017 and 159), and two of subsets (260 and 280) without other clues. Specifically, the extra information of solar elevation and solar azimuth serve as the input of sun position in the compared sat-nerf models. For subset 017 and 159 we take these extra sun information as sat-nerf input, and for subset 260 and 280 we take 45° as azimuth, and 90° elevation angle as the sun position placeholder for sat-nerf model.

Actually, for the sake of fairness, we shall keep all the input the same, like the process for subset 260 and 280. While here we make some concessions on subset 031 and 159 to validate the capability of SatensoRF for handling unknown sun position against Sat-NeRF with known sun positions.

Moreover, we also crop each subsets into overlapped multi-view imagery of size 512 × 512, to evaluate model performance on the same date. For every chosen original subsets of spacenet dataset, we crop the image indexed by 006 in the subset, namely 017-6, 033-6, 072-6 and 251-6. The cropping mode is shown in Fig. 8. And for single-date imagery we also use a subset 251-6 without extra sun information.

2) *Implemented Detail:* All the SatensoRF and its variants are implemented in PyTorch. For the baseline model we set  $L = 16$  levels of tensor decomposition for both the

<sup>2</sup><https://spacenet.ai/iarpa-multi-view-stereo-3d-mapping/>

<sup>3</sup><https://spacenet.ai/iarpa-multi-view-stereo-3d-mapping/>

TABLE I

THE DATASET COMPARISONS. ON THE MVS-DATASET FROM SPACENET, WE CREATE TWO TYPES OF DATASETS TO STUDY HOW SATENSORF ACTS ON SINGLE-DATE IMAGERY AND MULTI-DATE IMAGERY. # INDICATES THE COUNT OF NUMBER, AND THE MEASURE OF OVERLAP IS PIXEL. FOR SINGLE-DATE DATASET, WE CROP THE RGB IMAGE OF INDEX 006 INTO  $512 \times 512$  PATCHES, AND FOR MULTI-DATE DATASET, WE USE THE IMAGES OF THE ORIGINAL SIZES. ✓ INDICATES THE EXISTENCE OF THE ELEVATION AND AZIMUTH ANGLE OF SUN POSITION. ALT. BOUNDS INDICATES THE ALTITUDE RANGE OF THE AREA WITHIN THE IMAGE.

	subset	Resolution	#Rays	#Images	overlap[pix]	#train	#test	elevation	azimuth	Alt. bounds [m]
<b>single -date</b>	017-6	$512 \times 512$	262604	41	0.7813	18	2	✓	✓	[-23, 4]
	033-6	$512 \times 512$	262604	41	0.6836	18	2	✓	✓	[-27, 3]
	072-6	$512 \times 512$	262604	41	0.5859	18	2	✓	✓	[-22, 1]
	251-6	$512 \times 512$	262604	41	0.4883	18	2	✗	✗	[-28, 14]
<b>multi -date</b>	017	$2400 \times 2400$	5760000	10	-	9	1	✓	✓	[-23, 4]
	159	$2400 \times 2400$	5760000	11	-	10	1	✓	✓	[-29, 15]
	260	$2400 \times 2400$	5760000	16	-	14	2	✗	✗	[-30, 13]
	280	$2400 \times 2400$	5760000	24	-	21	3	✗	✗	[-29, 3]
<b>DFC</b>	004	$780 \times 765$	596276	11	-	9	2	✓	✓	[-24, 1]
	068	$795 \times 771$	613480	19	-	17	2	✓	✓	[-27, 30]
	214	$790 \times 761$	601292	24	-	21	3	✓	✓	[-29, 73]
	260	$794 \times 771$	612412	17	-	15	2	✓	✓	[-30, 13]

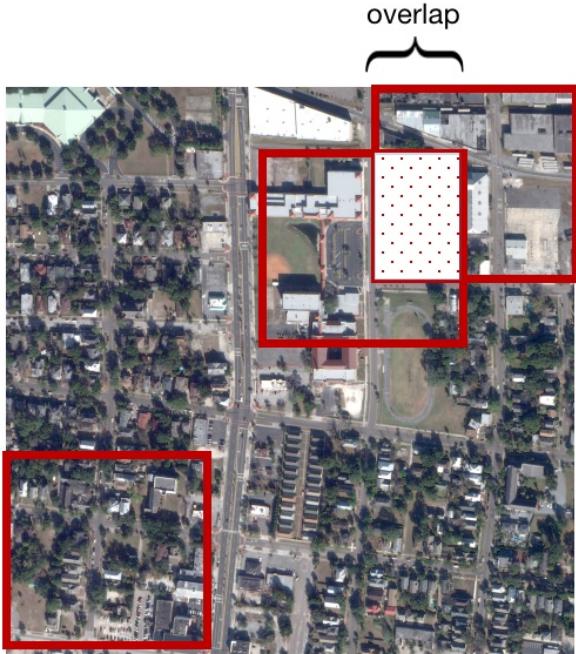


Fig. 8. Illustration of cropping single-date imagery with fixed-size of overlapping to support rational scene geometry for nerf-based method training.

appearance field  $\mathcal{T}_{ref}$  and the density field  $\mathcal{T}_\sigma$ ,  $L=1$  level for the sky color field  $\mathcal{T}_\sigma$ . The coefficient for the RGB loss  $\mathcal{L}_{rgb}$  and the total variation loss  $\mathcal{L}_{TV}$  are both 1.0. All the learnable parameters are optimized with Adam optimizer with a learning rate of 0.0002 for the tensor structure and 0.0001 for the MLPs. Both of the learning rates degrade log-linearly to 0.1 times to their values in the latest steps. The batchsize is set to 8192, and all the experiments are conducted on a single card of NVIDIA GeForce RTX 3090 GPU with 16 GB RAM.

### B. Efficiency Comparison

Firstly, we show that SatensoRF precedes other satellite neural radiance field in efficiency in both training and

inference phase. For multi-date dataset, we stop training when the validation PSNR starts to fall, and all of the time are taken average over different subsets of dataset. While for single-date dataset, since PSNR increases fast and we stop training after the validation PSNR reaches 36. Here we basically take the measures below to show the computation efficiency:

- 1) Training time  $T_{train}^{multi}(s)$  on multi-date dataset in seconds, until the validation PSNR starts to fall.
- 2) Training time  $T_{train}^{single}(s)$  on multi-date dataset in seconds, until the validation PSNR reaches 36dB.
- 3) Inference time  $T_{test}^{multi}(s)$  of computing a total image from the multi-date dataset in seconds.
- 4) Inference time  $T_{test}^{single}(s)$  of computing a total image from the single-date dataset in seconds.
- 5) Maximal setting of batchsize  $BS_{max}$  for the model to reach the maximum of the single-card GPU memory.
- 6) Number of parameters  $\#Params$ : number of the total parameters with  $M$  indicating a million.
- 7)  $FLOPs^{multi}$  (floating point operations) for multi-date imagery with  $B$  indicating a billion.
- 8)  $FLOPs^{single}$  (floating point operations) for single-date imagery with  $B$  indicating a billion.

The basic comparisons are listed in Table II. As can be seen, the full-size SatensoRF model has roughly a quarter of parameters of Sat-NeRF models, with the memory cost are far more less during training and testing, so the batchsize of SatensoRF can be set to nearly 8 to 16 times of the Sat-NeRF batchsize. Besides, the computation of SatensoRF are roughly one third of Sat-NeRF measured in FLOPs.

### C. Performance Comparison

The output of SatensoRF is twofold, a synthesized novel view and an altitude map, so we compare SatensoRF to the extant neural radiance fields for satellite with satellite cameras. The comparisons mainly include vanilla NeRF [7], S-NeRF [10], Sat-NeRF [11]. Since the experiments in [11] has validated the effectiveness of solar correction loss for

TABLE II  
EFFICIENCY COMPARISON OF SATENSORF AND THE STATE-OF-THE-ART METHODS.  $M$  IS THE ABBREVIATION FOR MILLION, AND  $B$  IS THE ABBREVIATION FOR BILLION.

	$T_{train}^{multi}(s)$	$T_{train}^{single}(s)$	$T_{test}^{multi}(s)$	$T_{test}^{single}(s)$	$BS_{max}$	#Params( $M$ )	$FLOPs^{multi}(B)$	$FLOPs^{single}(B)$
NeRF	26648.04	14481.13	131.47	17.23	1024	2.30	26462.45	1204.34
S-NeRF	28021.71	12280.29	136.76	18.02	1024	2.50	28768.67	1309.30
Sat-NeRF	37129.39	13826.52	141.27	19.31	1024	2.64	30293.36	1378.68
SatensoRF	<b>10129.81</b>	<b>2069.65</b>	<b>97.49</b>	<b>9.35</b>	<b>8096</b>	<b>0.79</b>	<b>416.21</b>	<b>416.21</b>

TABLE III  
NUMERICAL RESULTS COMPARISONS OF SATENSORF AND THE SOTA WORK ON MULTI-DATE IMAGES. POSTFIX -SC REFERS TO MODELS WITH SOLAR CORRECTION SUPERVISION IN SAT-NERF PAPER [11], AND -DS REFERS TO MODELS WITH DEPTH SUPERVISION.

Area Name	Test PSNR(dB)↑				Test SSIM↑				Test Altitude MAE(m)↓			
	031	159	260	280	031	159	260	280	031	159	260	280
NeRF	19.595	18.794	16.876	18.921	0.716	0.761	0.483	0.703	3.449	2.628	-	-
S-NeRF + SC	21.059	22.013	17.036	20.514	0.787	0.798	0.504	0.711	2.848	2.318	-	-
Sat-NeRF + SC	21.090	22.397	20.710	20.989	0.792	0.891	0.770	0.783	2.722	1.861	-	-
Sat-NeRF + DS + SC	21.617	23.812	21.094	21.322	0.801	0.901	0.792	0.813	<b>1.982</b>	<b>1.351</b>	-	-
SatensoRF + SC	24.629	26.695	22.706	24.578	0.921	0.933	0.918	0.939	2.943	2.081	-	-
SatensoRF + DS + SC	<b>24.771</b>	<b>27.013</b>	<b>23.151</b>	<b>24.866</b>	<b>0.934</b>	<b>0.929</b>	<b>0.952</b>	<b>0.931</b>	1.998	1.859	-	-

TABLE IV  
NUMERICAL RESULTS COMPARISONS OF SATENSORF AND THE SOTA WORKS ON THE SEEN CASES OF MULTI-DATE IMAGERY. POSTFIX -SC REFERS TO MODELS WITH SOLAR CORRECTION SUPERVISION IN SAT-NERF PAPER [11], AND -DS REFERS TO MODELS WITH DEPTH SUPERVISION.

Area Name	Train PSNR(dB)↑				Train SSIM↑				Train Altitude MAE(m)↓			
	031	159	260	280	031	159	260	280	031	159	260	280
NeRF	26.342	23.847	27.601	24.905	0.723	0.805	0.753	0.823	3.357	2.051	-	-
S-NeRF + SC	22.901	26.088	18.914	21.018	0.799	0.802	0.646	0.724	2.314	2.009	-	-
Sat-NeRF + SC	23.909	24.105	23.782	21.972	0.835	0.914	0.881	0.831	2.531	1.779	-	-
Sat-NeRF + DS + SC	24.101	24.597	24.710	22.849	0.821	0.909	0.804	0.821	<b>1.759</b>	<b>1.104</b>	-	-
SatensoRF + SC	29.639	32.705	30.988	29.618	0.932	<b>0.958</b>	0.924	<b>0.951</b>	2.815	2.019	-	-
SatensoRF + DS + SC	<b>29.872</b>	<b>33.382</b>	<b>31.210</b>	<b>29.905</b>	<b>0.956</b>	0.947	<b>0.962</b>	0.949	1.762	1.596	-	-

Sat-NeRF, we only take Sat-NeRF with solar-correction (SC) for comparison to yield better results, and the notations are clarified below:

- S-NeRF + SC: S-NeRF with solar correction.
- Sat-NeRF + SC: Sat-NeRF with solar correction.
- Sat-NeRF + DS: Sat-NeRF with depth supervision.
- Sat-NeRF + DS + SC: Sat-NeRF with solar correction and depth supervision.
- SatensoRF + SC: SatensoRF with Solar Correction.

1) *Multi-date Comparison*: Table III presents the multi-date comparison of SatensoRF and the SOTA neural radiance field works for multi-date imagery on test images unseen by the trained model. With shorter training time and inference time, SatensoRF achieves better results in novel view synthesis task with PSNR and SSIM as metrics, and the quality of the reconstructed DSM products equals to Sat-NeRF by and large. The MAE of subset 260 and 280 are omitted because there are water within the area, and we did not perform exact segmentations to find the accurate areas to compute DSM.

While Table IV presents the multi-date comparison for multi-date imagery on the training images. As shown in Table IV, SatensoRF converges faster to higher PSNR than all of the NeRF-based models, and the SSIM has increased more during shorter training time, synchronously with PSNR. The synthesized results with details of subsets 031, 159 are shown in Fig. 9, Fig. 9 respectively. The first and second rows show the different dates and the sun azimuth of the first 10 images synthesized by SatensoRF. The 3rd and the 5th row shows two levels of details in the synthesized images. The 6th row shows the estimated altitude map from the synthesized DSM, while the last row shows the ambient map ( $c_{amb}$ ) of different date.

Both of the two tables come to two main conclusion. Firstly, for multi-date images of original large size, SatensoRF has better performance than Sat-NeRF and previous methods. Secondly, from the results on subset 260 and 280, even though the input for SatensoRF has no solar information, the novel view synthesizing tasks still takes the lead in PSNR/SSIM.

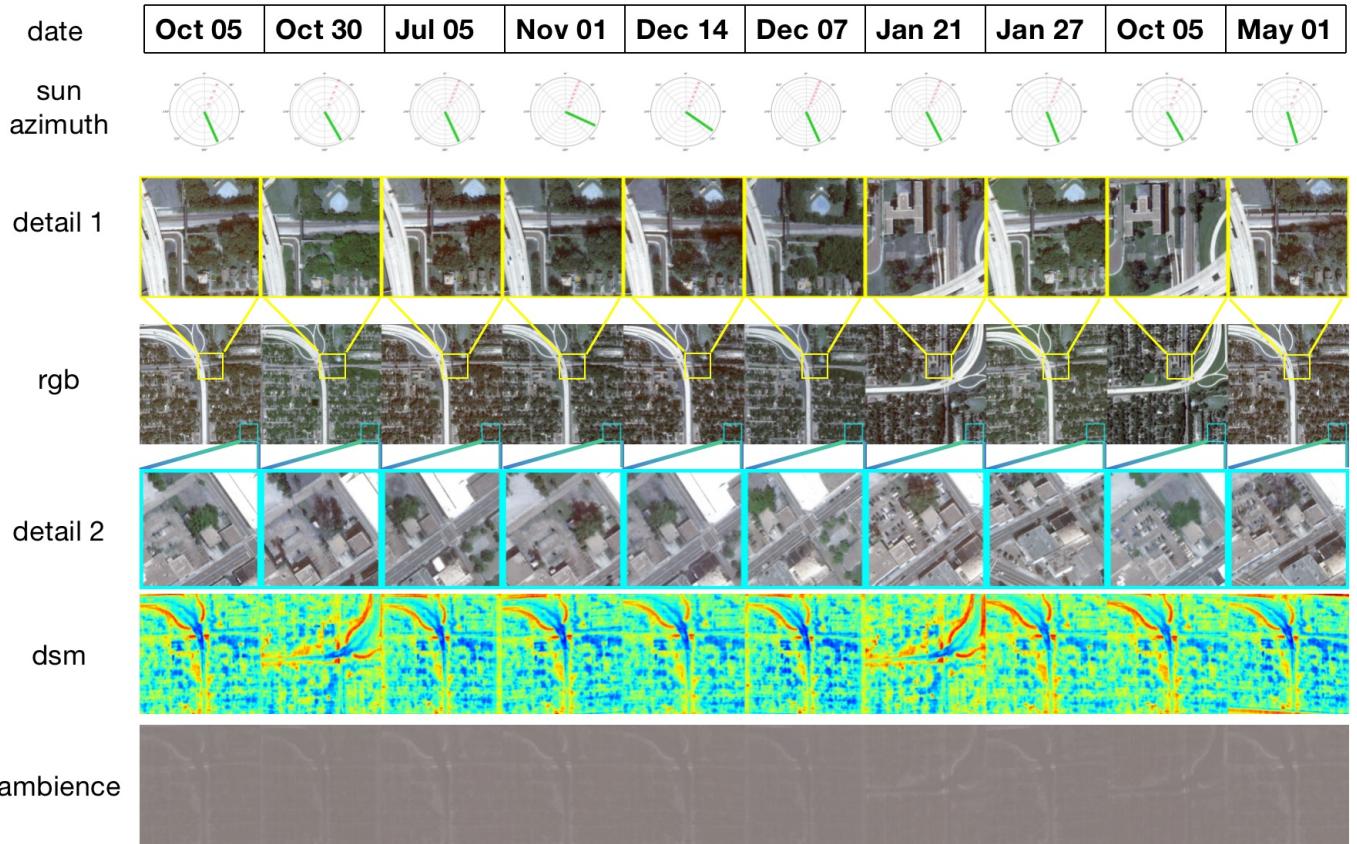


Fig. 9. Details of the synthesized images in subset 031, with details of street views of size  $500 \times 500$  and details of vegetations of size  $300 \times 300$ .

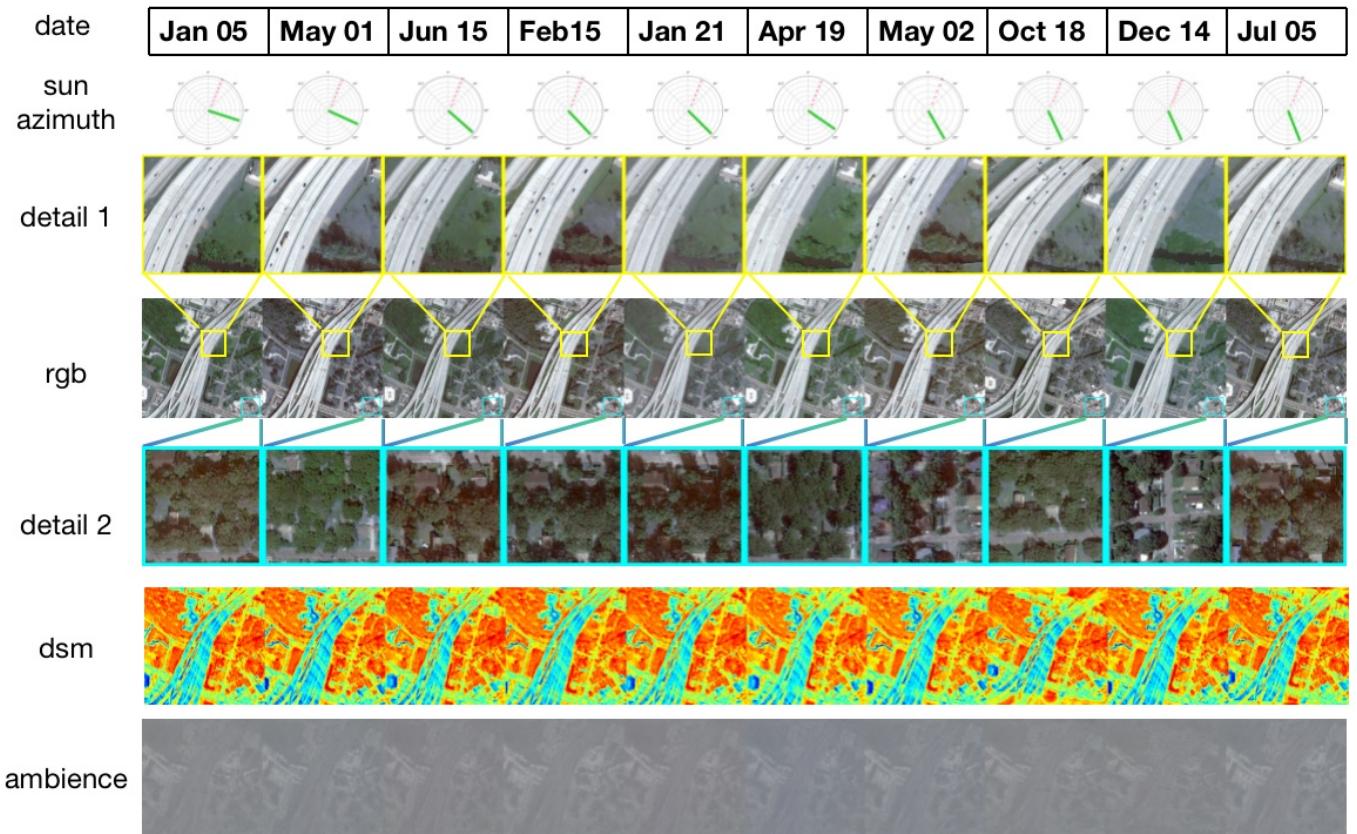


Fig. 10. Details of the synthesized images in subset 159, with details of asphalt flyover of size  $500 \times 500$  and details of street views of size  $300 \times 300$ .



Fig. 11. Details of the synthesized images in subset 260, with details of water body of size  $500 \times 500$  and details of vegetations of size  $300 \times 300$ .

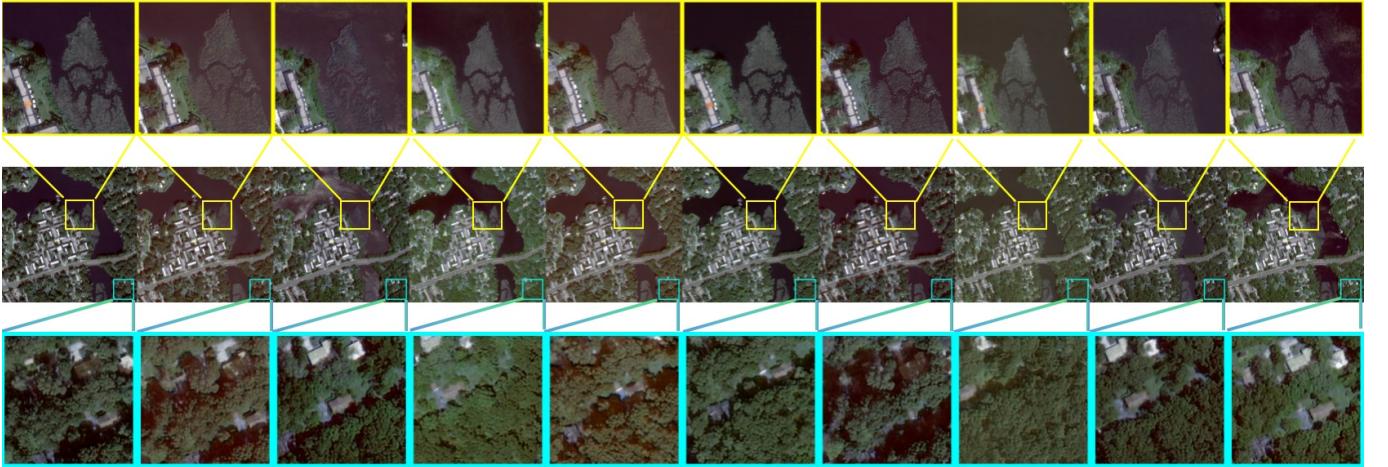


Fig. 12. Details of the synthesized images in subset 280, with details of water body of size  $500 \times 500$  and details of vegetations of size  $300 \times 300$ .

2) *Single-date Comparison*: In the same model zoo, for single-date imagery, Table V documents a comparative analysis between SatensoRF and state-of-the-art neural radiance field models. These evaluations were conducted on test and training images from a neighboring area. We stop trainings once the validation PSNR reaches a threshold of 37.0. Consequently, all models demonstrate comparable performance on the test set and differ a little on the training set. It is worth noting that, with the same lightning conditions, the additional modules of S-NeRF and Sat-NeRF does not promote synthesis quality than NeRF, while vanilla NeRF is faster. With similar performance, SatensoRF converges faster than all the NeRF-based models with the same level of performance.

#### D. Ablation Study I: Non Lambertian Assumption

This section aims to validate the non-Lambertian assumption by modeling the anisotropic reflectance field. The candidates for comparison are as follows:

- lamb CP + SH: The Lambertian version version of

SatenoRF without modeling the reflective light. The network has only the diffusion head (MLPs of 4 layers with 128 hidden neurons), and no other output from module A in Fig. 4. The final RGB output is rendered with spherical harmonics functions (SH), and the tensor factorization adopts CP (CANDECOMP/PARAFAC) decomposition.

- non lamb CP + asg: The non Lambertian version of SatensoRF modeling the reflective light with CP decomposition. The network is the same as module A in Fig. 4 and the final RGB is rendered with asg functions.
- lamb VM + SH: The same with lamb CP + SH but the tensor radiance fields are composed of VM (Vector-Matrix) decomposition.
- non lamb VM + asg: The SatensoRF itself.
- non lamb VM + SH: The same with non lamb VM + asg, but the final RGB is rendered with SH functions.

We evaluate the performance of these models on the test set and training set of subset 159. The quantitative results of

TABLE V

NUMERICAL RESULTS COMPARISONS OF SATENSORF AND THE SOTA WORK ON SINGLE DATE IMAGERY. POSTFIX *-SC* REFERS TO MODELS WITH SOLAR CORRECTION SUPERVISION IN SAT-NERF PAPER [11], AND *-DS* REFERS TO MODELS WITH DEPTH SUPERVISION.

Area Name	PSNR↑(test/train)				SSIM↑(test/train)			
	017-6	033-6	072-6	251-6	017-6	033-6	072-6	251-6
NeRF	37.012/42.561	36.984/43.784	36.971/43.659	36.969/42.998	0.901/0.923	0.982/0.988	0.913/0.968	0.948/0.952
S-NeRF + SC	36.925/45.032	35.779/43.291	36.216/42.145	35.968/40.328	0.896/0.972	0.878/0.984	0.915/0.983	0.904/0.936
Sat-NeRF + SC	36.546/43.592	37.004/44.672	36.882/44.105	36.829/43.045	0.935/0.972	0.921/0.986	0.905/0.928	0.931/0.979
SatenoRF + SC	<b>36.942/42.306</b>	<b>37.031/43.694</b>	<b>36.699/43.320</b>	<b>36.829/42.984</b>	<b>0.955/0.988</b>	<b>0.932/0.976</b>	<b>0.959/0.992</b>	<b>0.946/0.984</b>

TABLE VI

NUMERICAL RESULTS COMPARISONS OF SATENSORF WITH BRDF ASSUMPTIONS AND LAMBERTIAN ASSUMPTION.

	PSNR↑(test/train)	SSIM↑(test/train)	Altitude MAE(m)↓(test/train)	$T_{train}^{multi}(s)$	$T_{train}^{single}(s)$
lamb CP + SH	20.382/24.071	0.792/0.827	-	<b>8529.31</b>	<b>57.27</b>
non lamb CP + asg	21.518/24.599	0.781/0.847	-	10010.29	79.39
lamb VM + SH	25.473/30.074	0.872/0.901	3.104/2.945	10113.83	93.99
non lamb VM + SH	25.819/29.818	0.885/0.893	2.964/2.861	10138.73	95.85
non lamb VM + asg	<b>26.695/32.705</b>	<b>0.918/0.952</b>	<b>2.081/2.021</b>	10129.81	97.49

TABLE VII

NUMERICAL RESULTS COMPARISONS OF SATENSORF WITH TOTAL VARIATION REGULARIZATION AND L1 REGULARIZATION.

Regularization	PSNR↑(test/train)	SSIM↑(test/train)	Altitude MAE(m)↓(test/train)	$T_{train}^{multi}(s)$	$T_{train}^{single}(s)$
TV	<b>26.695/32.705</b>	<b>0.918/0.952</b>	<b>2.081/2.021</b>	10129.81	97.49
L1	26.314/31.059	0.857/0.886	-	<b>10096.53</b>	<b>97.29</b>
L1 + TV	26.518/31.622	0.872/0.894	-	10359.47	97.39

TABLE VIII

PERFORMANCE COMPARISON BY CHANGING THE TENSOR FIELD LEVEL  $L$ .

	PSNR(dB)↑	SSIM ↑	Altitude MAE(m) ↓	$T_{train}(s)$	$T_{test}(s)$	#Params(M)	FLOPs <sup>multi</sup> (B)
layer-2 $\mathcal{T}_\sigma$ -2	23.208	0.877	5.382	<b>8189.29</b>	<b>57.23</b>	0.24	129.52b
layer-2 $\mathcal{T}_\sigma$ -4	24.098	0.894	4.367	9101.337	70.405	0.25	132.74b
layer-4 $\mathcal{T}_\sigma$ -8	27.031	0.922	2.009	10452.052	105.31	0.86	450.57b
layer-2 $\mathcal{T}_\sigma$ -8(SatenoRF)	<b>26.695</b>	<b>0.918</b>	<b>2.081</b>	10129.813	97.495	0.79	416.21b

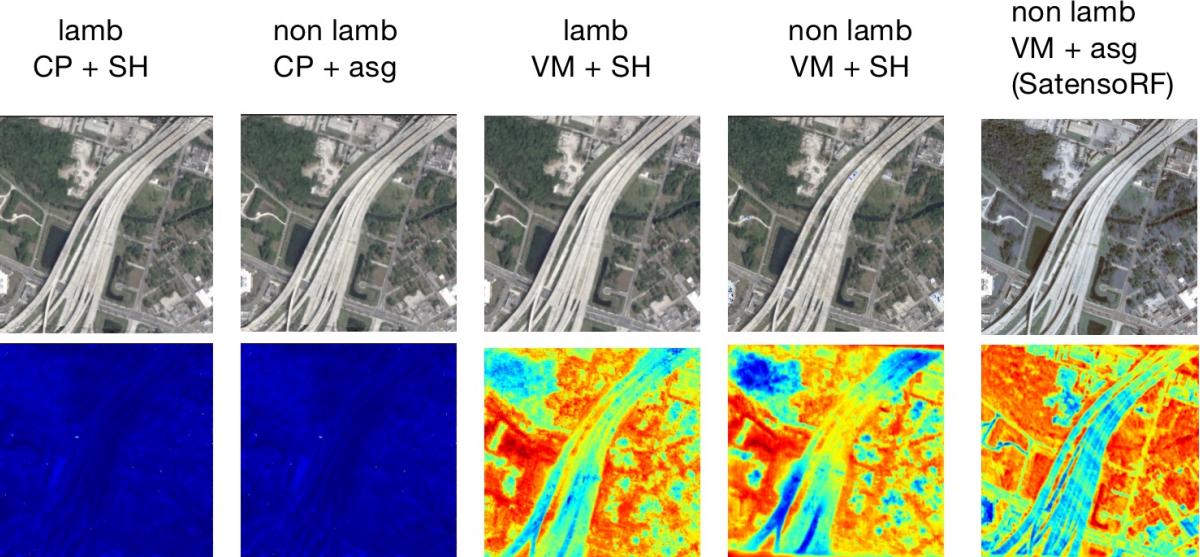


Fig. 13. The comparison of whether using Lambertian Assumption.

performance and the model are listed below in Table VI, and the visualizations are summarized in Fig. 13. From Fig. 13 we can see that the replacing VM decomposition with CP decomposition results in wrong volume density computation, and the wrong altitude map can not produce DSM product correspondingly so MAE is omitted in Table VI.

#### E. Ablation Study II: the Total Variation Loss

This section seeks to validate the rationale behind using the total variation loss, to mitigate the adverse effects of transient objects as noise on the volume density field  $\mathcal{T}_\sigma$ . In the default configuration of SatensoRF, the loss function specified in Equation 15 is employed.

The comparative analysis encompasses the following factors:

- TV: Training SatensoRF with  $\mathcal{L}_{TV}$ .
- L1: Training SatensoRF with  $\mathcal{L}_{l1}$ .
- L1 + TV: Training SatensoRF with both  $\mathcal{L}_{l1}$  and  $\mathcal{L}_{TV}$

where  $\mathcal{L}_{l1}(x) = \frac{1}{M} \sum_{i=1}^{M-1} |\mathcal{T}_\sigma|$ .

In Table VII, the validity of total variation loss  $\mathcal{L}_{TV}$  is significant since it preserves the right geometry of the area, while  $\mathcal{L}_{l1}(x)$  leads the model to the wrong volume density with nice rendering results. In Besides, Table VII shows that once  $\mathcal{L}_{l1}(x)$  is adopted, whether  $\mathcal{L}_{TV}$  involves can not guarantee the volume density estimation.

#### F. Controlled Experiments

SatenoRF is comprised of tensor composited fields and MLPs, and in this study, we delve into how the design of these two components influences the overall performance and efficiency of the model. For brevity, we denote model  $\mathcal{T}_\sigma - m$  to be the SatenoRF model with  $m$  level of tensor radiance fields; we denote model  $layer - n$  to be the SatenoRF model with MLP network with  $n$  layers.

Table VIII shows that both increasing the levels of tensor radiance field or the network layers enhance the performance, but the enhancement is limited while the time and computation costs is unworthy. Additionally, insufficient level of tensor radiance field causes poor result in 3D scene reconstruction.

## V. CONCLUSION

We propose a fast, level-of-detail tensorized neural radiance field called SatenoRF for large scale satellite novel view synthesis. The proposed radiance field efficiently speeds up the training and inference process, with better rendering quality. At the level of principles, SatenoRF firstly proposes to model the anisotropic light field rather than assuming the earth surface as Lambertian, to boost the radiance field performance. Moreover, SatenoRF are free from extra input such as solar directions. We conduct extensive experiments to evaluate SatenoRF on both performance and efficiency. SatenoRF shows better image rendering quality than SOTA

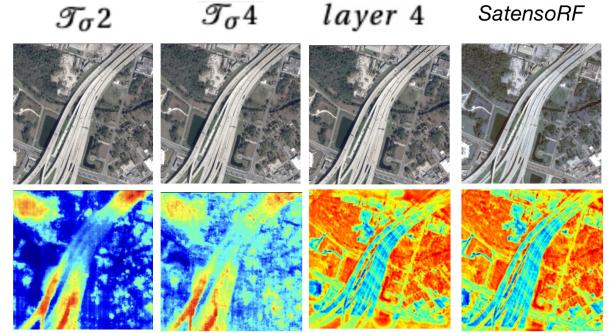


Fig. 14. The result comparison of the controlled experiments

Sat-NeRF series with much less computation, memory and training time on large satellite images. For future work, this paradigm can be extended with new dimension, such as time, to enable long-time ground surface observation.

#### ACKNOWLEDGMENT

This research is supported in part by the Special Funds for Creative Research Grant (No. 2022C61540).

#### REFERENCES

- [1] D. E. Shean, O. Alexandrov, Z. M. Moratto, B. E. Smith, I. R. Joughin, C. Porter, and P. Morin, "An automated, open-source pipeline for mass production of digital elevation models (dems) from very-high-resolution commercial stereo satellite imagery," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 116, pp. 101–117, 2016.
- [2] G. Facciolo, C. De Franchis, and E. Meinhardt-Llopis, "Automatic 3d reconstruction from multi-date satellite images," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017, pp. 57–66.
- [3] E. Rupnik, M. Pierrot-Deseilligny, and A. Delorme, "3d reconstruction from multi-view vhr-satellite images in micmac," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 139, pp. 201–211, 2018.
- [4] K. Zhang, N. Snavely, and J. Sun, "Leveraging vision reconstruction pipelines for satellite imagery," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [5] K. Gong and D. Fritsch, "Dsm generation from high resolution multi-view stereo satellite imagery," *Photogrammetric Engineering & Remote Sensing*, vol. 85, no. 5, pp. 379–387, 2019.
- [6] J. Gao, J. Liu, and S. Ji, "A general deep learning based framework for 3d reconstruction from multi-view stereo satellite images," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 195, pp. 446–461, 2023.
- [7] A. Tewari, O. Fried, J. Thies, V. Sitzmann, S. Lombardi, K. Sunkavalli, R. Martin-Brualla, T. Simon, J. Saragih, M. Nießner *et al.*, "State of the art on neural rendering," in *Computer Graphics Forum*, vol. 39, no. 2. Wiley Online Library, 2020, pp. 701–727.
- [8] D. Maggio, M. Abate, J. Shi, C. Mario, and L. Carlone, "Loc-nerf: Monte carlo localization using neural radiance fields," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 4018–4025.
- [9] H. Turki, D. Ramanan, and M. Satyanarayanan, "Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12922–12931.
- [10] D. Derksen and D. Izzo, "Shadow neural radiance fields for multi-view satellite photogrammetry," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 1152–1161.

- [11] R. Marí, G. Facciolo, and T. Ehret, "Sat-nerf: Learning multi-view satellite photogrammetry with transient objects and shadow modeling using rpc cameras," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 1311–1321.
- [12] R. Martin-Brualla, N. Radwan, M. S. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth, "Nerf in the wild: Neural radiance fields for unconstrained photo collections," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7210–7219.
- [13] W. Qin, J. R. Herman, and Z. Ahmad, "A fast, accurate algorithm to account for non-lambertian surface effects on toa radiance," *Journal of Geophysical Research: Atmospheres*, vol. 106, no. D19, pp. 22 671–22 684, 2001.
- [14] P. Hedman, P. P. Srinivasan, B. Mildenhall, J. T. Barron, and P. Debevec, "Baking neural radiance fields for real-time view synthesis," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5875–5884.
- [15] 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.
- [16] C. Reiser, R. Szeliski, D. Verbin, P. P. Srinivasan, B. Mildenhall, A. Geiger, J. T. Barron, and P. Hedman, "Merf: Memory-efficient radiance fields for real-time view synthesis in unbounded scenes," *arXiv preprint arXiv:2302.12249*, 2023.
- [17] M. Sun, D. Yang, D. Kou, Y. Jiang, W. Shan, Z. Yan, and L. Zhang, "Human 3d avatar modeling with implicit neural representation: A brief survey," *arXiv preprint arXiv:2306.03576*, 2023.
- [18] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Transactions on Graphics (ToG)*, vol. 41, no. 4, pp. 1–15, 2022.
- [19] 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.
- [20] C. Sun, M. Sun, and H.-T. Chen, "Improved direct voxel grid optimization for radiance fields reconstruction," *arXiv preprint arXiv:2206.05085*, 2022.
- [21] A. Chen, Z. Xu, A. Geiger, J. Yu, and H. Su, "Tensorf: Tensorial radiance fields," in *European Conference on Computer Vision*. Springer, 2022, pp. 333–350.
- [22] J. Tang, X. Chen, J. Wang, and G. Zeng, "Compressible-composable nerf via rank-residual decomposition," *Advances in Neural Information Processing Systems*, vol. 35, pp. 14 798–14 809, 2022.
- [23] K. Han and W. Xiang, "Multiscale tensor decomposition and rendering equation encoding for view synthesis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 4232–4241.
- [24] V. Rudnev, M. Elgarib, W. Smith, L. Liu, V. Golyanik, and C. Theobalt, "Nerf for outdoor scene relighting," in *European Conference on Computer Vision*. Springer, 2022, pp. 615–631.
- [25] H. Hirschmuller, "Stereo processing by semiglobal matching and mutual information," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 30, no. 2, pp. 328–341, 2007.
- [26] M. Michael, J. Salmen, J. Stallkamp, and M. Schlipsing, "Real-time stereo vision: Optimizing semi-global matching," in *2013 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2013, pp. 1197–1202.
- [27] K. Zhou, X. Meng, B. Cheng *et al.*, "Review of stereo matching algorithms based on deep learning," *Computational intelligence and neuroscience*, vol. 2020, 2020.
- [28] Y. Yao, Z. Luo, S. Li, T. Fang, and L. Quan, "Mvsnet: Depth inference for unstructured multi-view stereo," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 767–783.
- [29] X. Gu, Z. Fan, S. Zhu, Z. Dai, F. Tan, and P. Tan, "Cascade cost volume for high-resolution multi-view stereo and stereo matching," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2495–2504.
- [30] Y. Yao, Z. Luo, S. Li, T. Shen, T. Fang, and L. Quan, "Recurrent mvsnet for high-resolution multi-view stereo depth inference," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 5525–5534.
- [31] K. Luo, T. Guan, L. Ju, Y. Wang, Z. Chen, and Y. Luo, "Attention-aware multi-view stereo," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1590–1599.
- [32] S. Galliani, K. Lasinger, and K. Schindler, "Gipuma: Massively parallel multi-view stereo reconstruction," *Publikationen der Deutschen Gesellschaft für Photogrammetrie, Fernerkundung und Geoinformation e. V.*, vol. 25, no. 361-369, p. 2, 2016.
- [33] Q. Xu and W. Tao, "Planar prior assisted patchmatch multi-view stereo," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 12 516–12 523.
- [34] F. Wang, S. Galliani, C. Vogel, P. Speciale, and M. Pollefeys, "Patchmatchnet: Learned multi-view patchmatch stereo," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14 194–14 203.
- [35] K. Luo, T. Guan, L. Ju, H. Huang, and Y. Luo, "P-mvsnet: Learning patch-wise matching confidence aggregation for multi-view stereo," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 10 452–10 461.
- [36] A. Chen, Z. Xu, F. Zhao, X. Zhang, F. Xiang, J. Yu, and H. Su, "Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 14 124–14 133.
- [37] Y. Wei, S. Liu, Y. Rao, W. Zhao, J. Lu, and J. Zhou, "Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5610–5619.
- [38] K. Deng, A. Liu, J.-Y. Zhu, and D. Ramanan, "Depth-supervised nerf: Fewer views and faster training for free," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 882–12 891.
- [39] J. Chibane, A. Bansal, V. Lazova, and G. Pons-Moll, "Stereo radiance fields (srif): Learning view synthesis for sparse views of novel scenes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 7911–7920.
- [40] B. Kaya, S. Kumar, F. Sarno, V. Ferrari, and L. Van Gool, "Neural radiance fields approach to deep multi-view photometric stereo," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 1965–1977.
- [41] J. T. Kajiya, "The rendering equation," in *Proceedings of the 13th annual conference on Computer graphics and interactive techniques*, 1986, pp. 143–150.
- [42] M. Pharr, W. Jakob, and G. Humphreys, *Physically based rendering: From theory to implementation*. MIT Press, 2023.
- [43] H. Jin, I. Liu, P. Xu, X. Zhang, S. Han, S. Bi, X. Zhou, Z. Xu, and H. Su, "Tensorf: Tensorial inverse rendering," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 165–174.
- [44] K. Xu, W.-L. Sun, Z. Dong, D.-Y. Zhao, R.-D. Wu, and S.-M. Hu, "Anisotropic spherical gaussians," *ACM Transactions on Graphics (TOG)*, vol. 32, no. 6, pp. 1–11, 2013.
- [45] Y. Chang, L. Yan, H. Fang, and C. Luo, "Anisotropic spectral-spatial total variation model for multispectral remote sensing image destriping," *IEEE Transactions on Image Processing*, vol. 24, no. 6, pp. 1852–1866, 2015.

- [46] W. He, N. Yokoya, L. Yuan, and Q. Zhao, "Remote sensing image reconstruction using tensor ring completion and total variation," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 57, no. 11, pp. 8998–9009, 2019.
- [47] Z. Lin, "A review on low-rank models in data analysis," *Big Data and Information Analytics*, vol. 1, no. 2/3, pp. 139–161, 2016.
- [48] S. Gu and R. Timofte, "A brief review of image denoising algorithms and beyond," *Inpainting and Denoising Challenges*, pp. 1–21, 2019.