

Rational Polynomial Camera Model Warping for Deep Learning Based Satellite Multi-View Stereo Matching

Jian Gao* Jin Liu* Shunping Ji[†]

School of Remote Sensing and information Engineering, Wuhan University, China

{jian-gao, liujinwhu, jishunping}@whu.edu.cn

Abstract

Satellite multi-view stereo (MVS) imagery is particularly suited for large-scale Earth surface reconstruction. Differing from the perspective camera model (pin-hole model) that is commonly used for close-range and aerial cameras, the cubic rational polynomial camera (RPC) model is the mainstream model for push-broom linear-array satellite cameras. However, the homography warping used in the prevailing learning based MVS methods is only applicable to pin-hole cameras. In order to apply the SOTA learning based MVS technology to the satellite MVS task for large-scale Earth surface reconstruction, RPC warping should be considered. In this work, we propose, for the first time, a rigorous RPC warping module. The rational polynomial coefficients are recorded as a tensor, and the RPC warping is formulated as a series of tensor transformations. Based on the RPC warping, we propose the deep learning based satellite MVS (SatMVS) framework for large-scale and wide depth range Earth surface reconstruction. We also introduce a large-scale satellite image dataset consisting of 519 5120×5120 images, which we call the TLC SatMVS dataset. The satellite images were acquired from a three-line camera (TLC) that catches triple-view images simultaneously, forming a valuable supplement to the existing open-source WorldView-3 datasets with single-scanline images. Experiments show that the proposed RPC warping module and the SatMVS framework can achieve a superior reconstruction accuracy compared to the pin-hole fitting method and conventional MVS methods. Code and data are available at <https://github.com/WHU-GPCV/SatMVS>.

1. Introduction

Three-dimensional (3D) scene reconstruction from multi-view optical images is a hot topic in both the computer vision and remote sensing fields. Currently, there

are two mainstream camera models: the perspective camera model (i.e., the pin-hole model) and the cubic rational polynomial camera (RPC) model. The prevailing deep learning based multi-view stereo (MVS) methods [4, 5, 10, 20, 21, 31, 32, 33] have been designed and developed for multi-view images captured from pin-hole cameras. To date, the deep learning based MVS task for RPC-based images, i.e., optical satellite images, has not been tackled.

The study of satellite MVS task is very important as the multi-view satellite images captured from push-broom linear-array cameras with the RPC model are the main data source for large-scale 3D Earth surface reconstruction. In the recent deep learning based MVS methods [32, 5, 10], homography warping has been utilized to align images from different viewpoints to a reference view through a set of hypothetical fronto-parallel depth planes of the reference camera [32]. Analogously, developing a specific warping method for aligning multi-view satellite images to the reference image space could extend the deep learning based MVS methods to satellite images. The primary contribution of this work is that we propose, for the first time, a rigorous RPC warping module. Based on the RPC warping module, we then propose the SatMVS framework for deep learning based satellite MVS task. The SatMVS framework utilizes RPC warping to align the multi-view satellite images through the hypothetical height planes in the world coordinate system instead of the front-parallel depth planes. It features a multi-stage coarse-to-fine matching strategy on a feature pyramid to tolerate the wide search range, considering the great topographic relief of the Earth's surface and high-capacity satellite images. The SatMVS framework can also be used with most of the modern deep learning based MVS methods by replacing the perspective model with the RPC model and introducing the multi-stage strategy.

The other contribution is that we introduce a new open-source satellite MVS dataset (the TLC SatMVS dataset). We believe that this is the first open-source satellite dataset to be built from images collected by the push-broom three-line camera (TLC). The camera is mounted on the Ziyuan-3 (ZY-3) satellite. This dataset will be a beneficial com-

*Equally contributed.

[†]Corresponding author.

plement to the existing single linear-array WorldView-3 datasets such as MVS3D [3] and US3D [2]. The difference between a TLC and a single-scanline camera is that the former captures three images at the same scene simultaneously, and the latter has to shoot from different orbit positions at different times to form a stereo picture of the same scene.

In summary, we fill two gaps in this paper: 1) the lack of RPC warping module, which has hindered the application of the state-of-the-art MVS method to satellite images; and 2) the lack of large open-source TLC satellite MVS datasets, which has hampered the development of satellite MVS and large-scale Earth surface reconstruction.

2. Related Work

Datasets. Most public stereo datasets have been created for the reconstruction of close-range or small-scale scenes. Examples of the two-view stereo benchmark datasets are the Middlebury [26], KITTI[8], and SceneFlow[22] datasets. MVS benchmark datasets include the Middlebury-MVS[28], DTU[1], Tanks and Temples[18], and ETH3D[27] datasets. As a supplement to these close-range datasets, Liu and Ji[20] created the WHU MVS/Stereo dataset, which is the first large-scale multi-view aerial image dataset for the task of city-level Earth surface reconstruction. Multi-view satellite images are another important data source for Earth surface reconstruction. Currently, there are two public multi-view satellite datasets, both of which consist of WorldView-3 single linear array images. One is the MVS3D dataset[3], which provides 47 panchromatic images covering an area of approximately 100 km^2 near San Fernando, Argentina. However, the airborne LiDAR ground-truth data only cover a sub-area of about 20 km^2 , which is insufficient for the deep learning based MVS methods. The other is the US3D dataset[2], which is made up of 69 WorldView-3 images and covers a total area of 100 km^2 of two cities in the United States. The WorldView-3 single linear-array camera acquires MVS images on different dates. As a result, the temporal and seasonal changes in the MVS3D and US3D datasets cause obvious visual differences among the MVS images, leading to a negative effect on the MVS task.

3D Surface Reconstruction with Satellite Images. 3D Earth surface reconstruction from satellite imagery is mainly achieved through the traditional geometric methods, which can be roughly divided into two main types. The first type of method is based on the epipolar geometry of satellite images. An example of this is the RPC Stereo Processor (RSP)[24]. In this type, the stereo images are first rectified according to the RPC model[30], and a stereo matching algorithm such as the semi-global matching (SGM) stereo method [12] is then used to estimate the disparities. Finally, the disparity maps are converted into 3D points in the world coordinate system. The other type involves fitting a com-

plex RPC model into a pin-hole model in a small area, and then using the stereo/MVS pipeline for the reconstruction. Examples of this type are the satellite stereo pipeline (S2P) [6] for binocular stereo imagery and adapted COLMAP [34] for MVS imagery.

Cost Volume Construction. Cost volume construction is an essential part of the learning-based stereo/MVS methods. GC-Net [17] generates a 3D cost volume in stereo matching by concatenating the left and right feature maps along the disparity direction. [11] introduced a learned cost metric for MVS task by considering multi-patch similarity. SurfaceNet [16] and DeepMVS [13] pre-warp the multi-view images to the sweep planes in 3D space. Most of the recent state-of-the-art methods [4, 5, 10, 20, 21, 31, 32, 33] use differentiable homography warping to construct the cost volume or cost maps. Given a set of fronto-parallel planes of the reference camera at different depths, the images (or the extracted features) of each view are warped into the view of reference camera by a 3×3 homography matrix, and then the warped feature volumes are fused to a single cost volume. However, the homography matrix is derived from the pin-hole model, so that it is not suitable for multi-view satellite images with the RPC model. [34] proposed a solution by fitting the RPC model into a pin-hole model in small patches of a satellite image. However, this fitting strategy breaks the strict geometric relationship of the imaging model, and introduces inevitable errors, while also requiring a massive amount of preprocessing.

Deep Learning Based MVS Methods. The deep learning based MVS approaches can be classified into: 1) 3D convolution based methods, such as MVSNet [32], MVSNet++ [4], and P-MVSNet [21], which apply a series of 3D convolutions to regularize the cost volume; and 2) recurrent regularization based methods, such as R-MVSNet [33], which depth-wise process the cost maps of different depths. The former methods are intuitive, but require much more GPU capacity. All of these methods were developed for use with natural images. RED-Net [20] extends the convolutional gated recurrent unit (ConvGRU) based regularization method [33] for aerial MVS task. Methods based on building cost volumes at multiple stages have also been recently introduced, e.g., CasMVSNet [32], CVP-MVSNet [31], and UCS-Net[5], in which a coarse-to-fine pyramid matching structure is applied. These multi-stage methods narrow the depth search range of the current stage to construct a thin cost volume with lower memory and a higher depth-wise sampling rate, which is conducive to large-scale scene reconstruction. However, it is difficult to apply the state-of-the-art deep learning based MVS methods to satellite imagery with a complex RPC model. In this paper, we attempt to fill this gap by proposing a general deep learning based MVS framework for satellite images (SatMVS) with a novel rigorous RPC warping module.

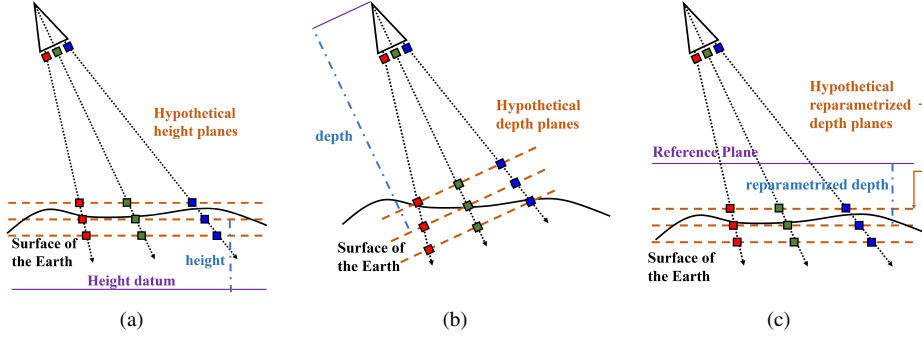


Figure 1: The hypothetical (a) height planes in our RPC warping, (b) depth planes, and (c) reparametrized depth planes.

3. RPC Warping

3.1. Rational Polynomial Camera Model

The RPC model is the most widely used camera model in high-resolution satellite images. It connects the image points and corresponding world coordinate points with cubic rational polynomial coefficients. In Eq.(1), a 3D point in the normalized world coordinates (*latitude, longitude, height*), denoted as (lat_n, lon_n, hei_n), is transformed to image space to obtain the normalized image coordinates ($samp_n, line_n$), which correspond to the along-array and along-track directions, respectively. Eq.(2) is the reverse version. P^{fwd} and P^{inv} are both cubic polynomials, as shown in Eq.(3), where the summation of integers m_1, m_2 and m_3 is no more than 3.

$$\begin{cases} samp_n = \frac{P_1^{fwd}(lat_n, lon_n, hei_n)}{P_2^{fwd}(lat_n, lon_n, hei_n)} \\ line_n = \frac{P_3^{fwd}(lat_n, lon_n, hei_n)}{P_4^{fwd}(lat_n, lon_n, hei_n)} \end{cases} \quad (1)$$

$$\begin{cases} lat_n = \frac{P_1^{inv}(samp_n, line_n, hei_n)}{P_2^{inv}(samp_n, line_n, hei_n)} \\ lon_n = \frac{P_3^{inv}(samp_n, line_n, hei_n)}{P_4^{inv}(samp_n, line_n, hei_n)} \end{cases} \quad (2)$$

$$P(X, Y, Z) = \sum_{i=0}^{m_1} \sum_{j=0}^{m_2} \sum_{k=0}^{m_3} c_{ijk} \cdot X^i \cdot Y^j \cdot Z^k \quad (3)$$

Clearly, the RPC model is a general geometric model, instead of a physical camera model; however, it has been proved that the RPC model can achieve a very high accuracy comparable to the rigorous sensor model (RSM) [29, 9], guaranteeing its successful application to all of the high-resolution optical satellite images.

3.2. Differentiable RPC Warping

In the pin-hole camera model, the geometric transformation between the corresponding pixels of stereo images can be formulated as a simple 3×3 homography matrix by given depths. Almost all of the state-of-the-art MVS methods [4, 5, 10, 20, 21, 31, 32, 33] warp the source views to

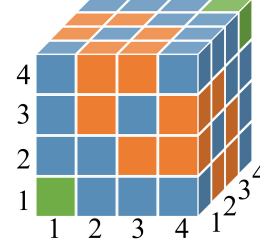


Figure 2: Coefficient tensor T . When i, j , and k are equal, $T_{ijk} = a_i a_j a_k$ (green blocks); when there are only two of i, j , and k equal, $T_{ijk} = a_i a_j a_k / 3$ (blue blocks); when none of i, j , and k are equal, $T_{ijk} = a_i a_j a_k / 6$ (orange blocks).

the reference view through a homography matrix and a series of hypothetical fronto-parallel planes of the reference view. In contrast, the RPC model is much more complicated and its warping cannot be formulated with only a matrix. In this paper, we propose a rigorous and efficient RPC warping module, which is fundamentally a high-order tensor transformation. It projects images (or extracted features) from different views to the reference view through a set of hypothetical height planes in the world coordinate system, instead of the fronto-parallel planes of a reference view, because there are no explicit physical parameters in RPC model to define what is the front of a camera. Fig. 1 shows the difference between the hypothetical height planes we used, the hypothetical depth planes, and reparametrized depth planes used in [Zhang2019LeveragingVR].

We start by transforming the ternary cubic polynomial in Eq.(3) to a quaternion cubic homogeneous polynomial in Eq.(3) to a quaternion cubic homogeneous polynomial $f(x_1, x_2, x_3, x_4) = \sum (a_i a_j a_k) \cdot x_i x_j x_k$, where $x_1 \equiv 1$ and $a_i a_j a_k (i, j, k \in \{1, 2, 3, 4\})$. By setting $x_2 = lon_n$, $x_3 = lat_n$, and $x_4 = hei_n$, f becomes the numerator or denominator of the RPC forward form, and that of the inverse form by setting $x_2 = line_n$, $x_3 = samp_n$, $x_4 = hei_n$. Please note that we use Einstein's summation convention in all of the formulas in this section. The four variables x_1, x_2, x_3, x_4 are then expressed as a tensor \mathbf{X} with a rank of 1: $\mathbf{X} = (x_1, x_2, x_3, x_4)^T$, and the polynomial coefficients are expressed as a tensor \mathbf{T} with a rank of 3 and a shape of

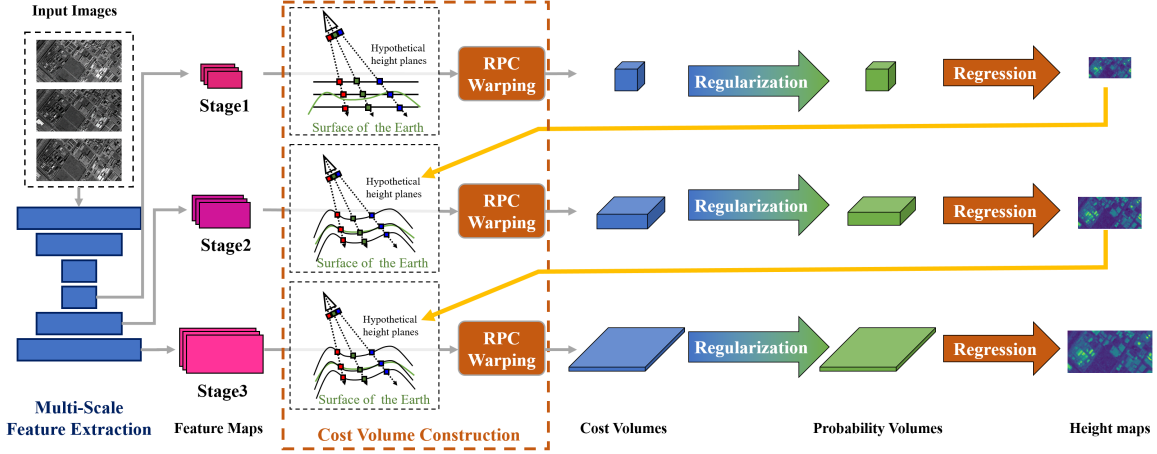


Figure 3: The structure of the SatMVS framework.

$4 \times 4 \times 4$ (see Fig. 2). Through the tensor contraction operation, which calculates the sum over all the possible values of the repeated indices of a set of tensors [23], the numerator and denominator of the RPC model can be expressed in tensor form:

$$f(\mathbf{X}) = \mathbf{T}_{ijk} \mathbf{X}_i \mathbf{X}_j \mathbf{X}_k \quad (4)$$

The relationship between the elements in \mathbf{T}_{ijk} and $a_i a_j a_k$ is: when i, j , and k are equal, $\mathbf{T}_{ijk} = a_i a_j a_k$; when there are only two of i, j , and k equal, $\mathbf{T}_{ijk} = a_i a_j a_k / 3$; when none of i, j , and k are equal, $\mathbf{T}_{ijk} = a_i a_j a_k / 6$. Finally, the left side of Eq.(1) or (2) can be easily obtained by dividing the numerator and denominator in an element-wise way. We can then extend this calculation to adapt to both the batch operation and the RPC transformation with a set of points:

$$f^{(bm)}(\mathbf{X}) = \mathbf{T}_{ijk}^{(b)} \mathbf{X}_i^{(bm)} \mathbf{X}_j^{(bm)} \mathbf{X}_k^{(bm)} \quad (5)$$

where $X^{(bm)}$ represents the m -th point in the b -th batch and $T^{(b)}$ represents the coefficient tensor in the b -th batch. The elements in parentheses do not participate in the summation. Through element-wise division, the RPC warping of all the points in a batch can be calculated in one shot.

Specifically, for a point $(samp_s, line_s)$ in the source image and a given plane hei , we obtain the normalization form $(samp_{s,n}, line_{s,n}, hei_n)$ according to the available normalization parameters, which are also a part of the RPC parameters, and construct tensor $\mathbf{X}_s = (1, samp_{s,n}, line_{s,n}, hei_n)$, and then warp it onto the reference view through:

$$\begin{cases} samp_{r,n} = \frac{f_1^{fwd}(\mathbf{X}_0)}{f_2^{fwd}(\mathbf{X}_0)} \\ line_{r,n} = \frac{f_3^{fwd}(\mathbf{X}_0)}{f_4^{fwd}(\mathbf{X}_0)} \end{cases} \quad (6)$$

where $\mathbf{X}_0 = \left(1, \frac{f_1^{inv}(\mathbf{X}_s)}{f_2^{inv}(\mathbf{X}_s)}, \frac{f_3^{inv}(\mathbf{X}_s)}{f_4^{inv}(\mathbf{X}_s)}, hei_n\right)$, and $\frac{(\cdot)}{(\cdot)}$ represents element-wise division.

After anti-normalization, the point $(samp_s, line_s)$ is warped to the corresponding point $(samp_r, line_r)$ with differentiable resampling, e.g., bilinear interpolation, to complete the RPC warping.

4. The Learning Based SatMVS Framework

We propose a satellite MVS deep learning framework imbedded with RPC warping, which we call the SatMVS framework. In addition to the necessary part of RPC warping, we consider that coarse-to-fine multi-stage processing is also necessary in such a framework to predict a wide range of Earth surface elevation. The other parts, such as the feature extraction, cost map regularization, and regression module, can be borrowed from the state-of-the-art methods [10, 20, 31]. Finally, we obtain the complete SatMVS framework, as shown in Fig. 3.

4.1. Multi-scale Feature Extraction

The current multi-scale MVS methods utilize popular feature extractors such as a feature pyramid network [19] or UNet [25]. All of these can be used in the SatMVS framework. For the experiments, we adopted a weight-shared multi-scale feature extractor [5], which consists of an encoder and a decoder with skip connections, to extract the features. The module outputs a three-scale feature pyramid with the size of $\{1/16, 1/4, 1\}$ of the input image size, and the numbers of channels are 32, 16, and 8, respectively.

4.2. Multi-Stage Cost Volume Construction with RPC Warping

There may be huge elevation differences in the hundreds of kilometers of landscape covered by a satellite image. We suppose that the maximum height difference is $2km$. If the height interval is set to $2.5m$, which is roughly the pixel resolution of ZY-3 imagery, 800 hypothetical planes can cover the entire elevation range with a pixel-level accu-

racy. In contrast, the depth search range of the close-range datasets is relatively small (e.g., 128 or 256 planes in the DTU dataset). To reduce the high demand for GPU memory and increase the learning speed, multi-scale learning from the feature pyramid is applied.

At the first stage ($i = 1$), the search scope should cover the range between the maximum and minimum heights of the covered area. This information can be found in the available RPC parameters. Alternatively, several open-source global digital elevation model (DEM) products, such as the Shuttle Radar Topography Mission (SRTM) DEM [15] or the Advanced Spaceborne Thermal Emission and Reflection Radiometer (ASTER) global digital elevation model (GDEM) [7] can be used as information sources. The number of hypothetical planes is fixed in the first stage (e.g., 64). The height interval is then determined by dividing the height difference by the number of hypothetical planes.

For the stage i ($i \geq 2$), the interval and the number of hypothetical planes are empirically fixed, and the hypothetical planes are centered at the reference height of the last stage. The feature maps from the source views are warped to the reference view through RPC warping to form multiple feature volumes, which are then fused to the cost volume using a variance-based operation [32].

4.3. Regularization

The recent MVS networks use either a series of 3D convolutions [32] or 2D convolutional GRUs [33] to regularize the cost volume. Both can be used in the SatMVS framework. In practice, we use the recurrent encoder-decoder structure of RED-Net [20] to regularize the cost maps that are constructed at each stage. Compared to 3D convolutions, the RED structure sequentially regularizes the cost maps along the elevation direction to achieve a high efficiency and low memory cost, which significantly benefits the processing of high-capacity satellite images with a wide height search range.

4.4. Height Inference

After regularization, a *soft argmin* operation is applied along the height direction for sub-pixel estimation. At the training stage, the pyramid network outputs height maps at three resolutions. Similar to the current works [2, 10], the total loss is defined as the weighted sum of the three-stage L1 loss, in which the weights are {0.5, 1, 2}, respectively.

4.5. Pipeline

We designed a complete pipeline for the RPC-based networks to reconstruct the final DSM, including 1) image preprocessing; 2) MVS inference; and 3) DSM generation. In the preprocessing, the study area in the world coordinate system (usually the WGS-84 coordinate system) is divided into regular blocks, and each block is projected onto

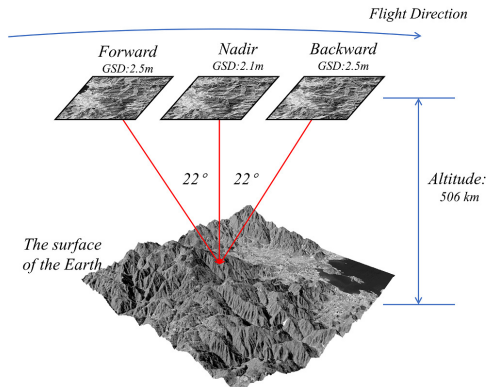


Figure 4: The three-line camera (TLC) of the ZY-3 satellite.

the MVS satellite images with the maximum and minimum elevation planes in this block respectively. The minimum bounding rectangles of the two projection regions in different views are calculated, extended to a uniform size, and cropped as the input of the network. This process ensures sufficient overlap between the cropped image patches of multiple views. In addition, gamma correction and linear enhancement can be used to improve the contrast.

In the MVS inference, the cropped images are fed into the proposed SatMVS framework. Each view image is treated as the reference image to infer the height map in turn.

In the DSM generation, the geometric consistency [32] is used to filter out outliers in the inferred height maps of the different views. If the distance between the reprojected point from the j -th view and the original point in the reference view is less than 1 pixel, the estimated result is considered to be geometrically consistent and valid. The valid matching results are then transformed to the Universal Transverse Mercator (UTM) coordinate system and resampled to a regularized DSM.

We also propose the homography warping version of the SatMVS pipeline. For the homography warping based networks, which cannot be directly used for satellite imagery, the RPC model of each cropped image is fitted to the pin-hole model according to [34] in the image preprocessing part. In the parts of MVS inference and DSM generation, the fitting model is consistently used instead of the rigorous RPC model.

5. The Satellite MVS Dataset

5.1. Data Source

This section describes the satellite MVS dataset we built, which is named the TLC SatMVS dataset. The triple-view images were collected from the TLC camera mounted on the Ziyuan-3 (ZY-3) satellite. The ground resolution of the nadir and the two side-looking images is 2.1 m and 2.5 m, respectively (see Fig. 4). As a professional satellite for surveying and 3D mapping, the ZY-3 satellite ac-

Image size(pixels)	768 ²	4608 ²	9216 ²	13824 ²	18432 ²	23040 ²
Min fitting error(pixels)	0.00015	0.00010	0.00181	0.00152	0.00106	0.00257
Max fitting error(pixels)	0.16204	1.09168	2.42714	3.89724	5.35514	6.62046

Table 1: The error of fitting an RPC model with a pin-hole model increases with the size of the image patch.

cesses the same scene at almost the same time, without the impact of illumination and seasonal changes, which differs from the WorldView series with single linear-array cameras. The RPC parameters have been calibrated in advance to achieve a sub-pixel reprojection accuracy. The ground-truth DSMs were prepared from both high-accuracy LiDAR observations and ground control point (GCP)-supported photogrammetric software [14]. The DSM was stored as a 5-m resolution regular grid under the WGS-84 geodetic coordinate system and the UTM projection coordinate system.

5.2. The TLC SatMVS Dataset

We built two versions of the TLC SatMVS dataset. The first version is a collection of large-size satellite images, and the second is a ready-made version for the training and testing of a learning method with mainstream GPU capacity.

In the first version (see Fig. 5), there are 173 sets of images (we call one triple-view image a set), with 127 sets separated for training and the rest for testing. Each set contains 16-bit panchromatic triple-view images with a 5120×5120 pixel size, the RPC parameters, and the ground-truth DSM, covering approximately 125 km^2 . The overlap rate of the triple-view images is more than 95%.

The second version (ready-made version) contains cropped patches of the images (see Fig. 6), the RPC parameters of each patch, and the corresponding height maps, which were obtained by projecting the DSMs to the images with the RPC parameters. The height map is theoretically equivalent to the depth map in a close-range MVS dataset, but stores the height instead of the depth information of the corresponding pixel in the image. Specifically, each 5120×5120 image is cropped into 768×384 patches with an overlap rate of 5% in both the horizontal and vertical direction. There are 5011 sets for training in total.

We also provide an accessory version, where we fitted each patch from the RPC projection into pin-hole projection according to [34] under the UTM coordinate system. The depth maps were obtained by projecting the ground truth to the patches with the fitted pin-hole camera parameters.

6. Experiments

6.1. Fitting the RPC Model with the Pin-hole Model

Before this work, the complex RPC model was fitted with a pin-hole model in the conventional satellite MVS [34]. The fitting error is unavoidable, and increases with the size of the image patch. As shown in Table 1, the maximum fitting error is small in the patch with a size of

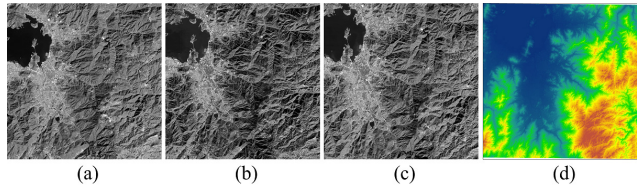


Figure 5: Examples of the collection of large-size satellite images ((a), (b), and (c) are the backward, forward, and nadir view images) and the corresponding ground truth (d) from the TLC SatMVS dataset.

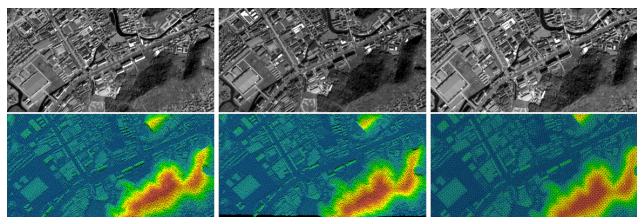


Figure 6: The ready-made version of the TLC SatMVS dataset for the training. The images from left to right are the image patches of the backward, forward, and nadir views, with the corresponding height maps below.

768×768 pixels, but it can reach more than 6 pixels at a size of 23040×23040 pixels, which is approximately the size of a ZY-3 TLC image or a WorldView-3 image. Fig. 7 also shows the error distributions in the XY-planes. Note that there are no geometric errors in the rigorous RPC model.

6.2. Model Evaluation

Implementation Details. We evaluated the proposed SatMVS framework with imbedded RPC warping on the TLC SatMVS dataset. The framework was implemented in PyTorch and trained on a single NVIDIA TITAN RTX GPU (24GB). Different MVS architectures, including RED-Net, CasMVSNet, and UCS-Net, were integrated into the proposed framework. The hyper-parameters followed the same settings in all of the experiments: in the training phase, the batch size was set to 1, and RMSprop was selected as the optimizer. All the networks were trained for 35 epochs with an initial learning rate of 0.001, and were downscaled by a factor of 2 after the 10th epoch.

Three-stage hierarchical matching was adopted to infer the coarse-to-fine height or depth maps. For TLC images, the view number of the input images N is fixed to 3. The numbers of hypothetical height planes were set to $\{64, 32, 8\}$ in the three stages, and the corresponding interval was $\{(d_{max} - d_{min})/64, 5m, 2.5m\}$, except for the UCS-Net implementation, where its own adaptive interval determina-

Methods	MAM(m)	RMSM(m)	<2.5m(%)	<7.5m(%)	Comp.(%)	Runtime
adapted COLMAP[34]	2.227	5.291	73.35	96.00	79.10	77min27s
RED-Net[20]*	2.171	4.514	74.13	95.91	81.82	9min15s
CasMVSNNet[10]*	2.031	4.351	77.39	96.53	82.33	4min02s
UCS-Net[5]*	2.039	4.084	76.40	96.66	82.08	3min47s
SatMVS(RED-Net)	1.945	4.070	77.93	96.59	82.29	13min52s
SatMVS(CasMVSNNet)	2.020	3.841	76.79	96.73	81.54	12min20s
SatMVS(UCS-Net)	2.026	3.921	77.01	96.54	82.21	13min17s

Table 2: Quantitative results of the different MVS methods on the TLC SatMVS dataset. The proposed SatMVS with RPC warping implements three different learning-based MVS methods for height inference. Recent deep MVS methods imbedded with fitted homography warping are marked with *.)

Methods	MAE(m)	RMSE(m)	<2.5m(%)	<7.5m(%)	Comp.(%)	Runtime
RED-Net (2048 × 1472)	2.171	4.515	74.13	95.91	81.82	9min10s
RED-Net (5120 × 5120)	2.517 (+0.346)	4.873 (+0.358)	66.42 (-7.71)	95.53 (-0.38)	81.44 (-0.38)	4min17s
SatMVS(RED-Net) (2048 × 1472)	1.945	4.071	77.93	96.59	82.29	13min12s
SatMVS(RED-Net) (5120 × 5120)	1.946 (+0.001)	4.224 (+0.153)	77.88 (-0.05)	96.54 (-0.05)	82.35 (+0.06)	5min52s

Table 3: Quantitative results of the SatMVS(RED-Net) and the RED-Net(with fitted pinhole model) on the TLC SatMVS dataset with different sizes.

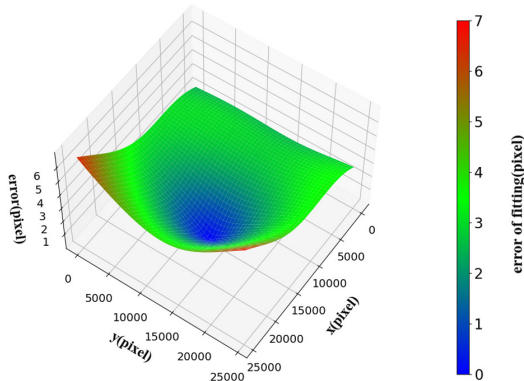


Figure 7: Error distribution of the RPC-to-pinhole model fitting. The X and Y coordinates represent the column and row coordinates of the image patches, and the Z axis represents the fitting error.

tion strategy was applied.

We also performed experiments on the MVS pipeline, with the fitted homography warping module used to replace the RPC warping module. Please note that double-precision floating point is used here to get rid of the numerical precision issue. As for adapted COLMAP [34], since it is itself a complete conventional satellite MVS pipeline, we directly used its own framework for the reconstruction.

We adopt four commonly used metrics to evaluate the quality of the final DSM: 1) the mean absolute error (MAE), i.e., the average of the L1 distance over all the grid units between the ground truth and the estimated DSM; 2) the root-mean-square error (RMSE), i.e., the standard deviation of the residuals between the ground truth and the estimation;

3) the percentage of grid units with an L1 distance error below the thresholds of 2.5 m (approximately equivariant to the ground sample distance (GSD)) and 7.5 m (<2.5m and <7.5m); and 4) the completeness (Comp.), i.e., the percentage of grid units with valid height values in the final DSM.

Evaluation on the TLC SatMVS Dataset. As there are currently no deep learning based MVS methods for satellite images, we only compared the proposed SatMVS framework with a recent conventional method, i.e., adapted COLMAP [34], which uses the pin-hole camera model to fit the RPC model. In addition, we imbed the recent learning based methods [5,10,20] with the fitted homography warping into our pipeline to generate DSMs for comparison. In contrast, our SatMVS applies the rigorous RPC warping with different MVS regularization methods. In the pipeline the images of the test set are cropped into patches with a size of approximately 2048×1472 pixels for inference. The results are listed in Table 2.

Several conclusions can be drawn from Table 2. Firstly, all of the learning based MVS methods perform better than the conventional adapted COLMAP method. The advantage is apparent in the RMSE metric, which indicates that the learning methods have a lower variance of inference. Fig. 8 shows samples of the reconstructed DSMs produced by RED-Net(with fitted homography warping), SatMVS(RED-Net), and adapted COLMAP. The regions that failed to match are colored in white, most of which suffer from challenging scenes (occlusions, clouds, shadows, and textureless water surfaces). Compared with adapted COLMAP, the DSM results of the learning methods are

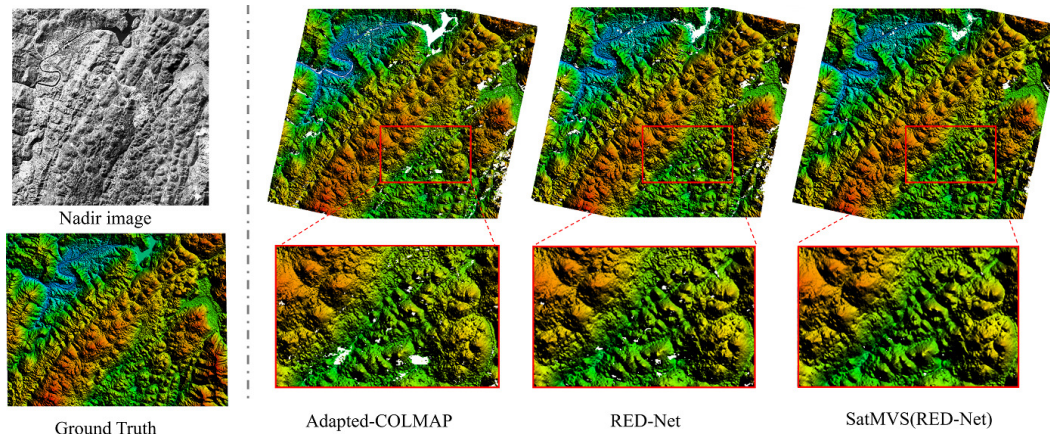


Figure 8: Results of RED-Net(homo), RED-Net(RPC), and adapted COLMAP[34].

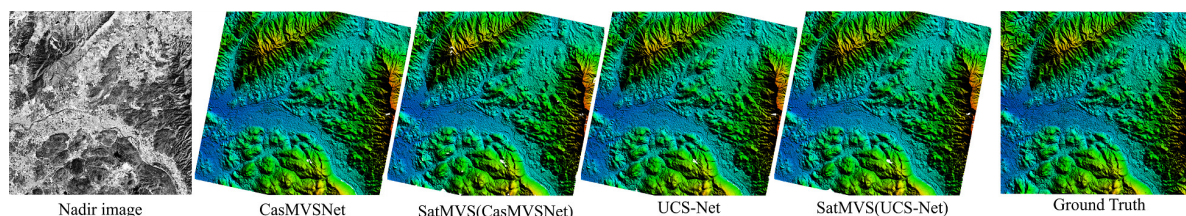


Figure 9: The qualitative results of CasMVSNet and UCS-Net with RPC warping and homography warping.

more complete, especially the RPC warping versions. Secondly, the performance of using the rigorous RPC model is only slightly better than using the fitted model. The reason is that the fitted model can also reach a sub-pixel accuracy at the size of 2048×1472 pixels, as shown in Table 1. Thirdly, the inference of the deep learning based methods is much faster than that of the traditional adapted COLMAP method. However, the speed of RPC warping is lower than that of homography warping as the latter is only built on a simple 3×3 matrix multiplication operation.

Evaluation on larger images. We process the full-size images (5120×5120 pixels) with RED-Net and our SatMVS(RED-Net) on the NVIDIA RTX A6000 GPU(48GB). The results are shown as Table 3. When processing 5120×5120 images, it is observed that, (1) the performance of RED-Net severely decreases due to the increase of fitting error, but our RPC warping method remains stable; (2) SatMVS significantly outperforms RED-Net, which confirms the effectiveness and advantage of the proposed method; (3) the efficiency of both methods is significantly improved and reaches a comparable level. In addition to the high efficiency, using larger-capacity images have significantly simplified the process because cropping images into highly overlapped small patches is difficult when DSM is unavailable and processing the small patches will cause repeated and redundant matching in overlapped regions. It should be mentioned that here we only use RED-Net as the memory consumption of using CasMVSNet and UCS-Net is unaffordable.

7. Discussion

We did not perform any more experiments on other datasets. The reason for this is that there are very few datasets that are suitable for the satellite MVS problem. The existing MVS3D dataset [3] is very small, and was designed for testing non-learning based methods. The US3D dataset[2] is aimed at the joint reconstruction and segmentation task, where the stereo images of different views vary dramatically in illumination and seasonal changes. We failed to train an effective deep learning model with either of these datasets. The presented TLC SatMVS dataset greatly relieve this situation and will contribute to the study of deep learning based satellite MVS method. Furthermore, the TLC camera is more suitable for Earth surface reconstruction than the single linear array cameras such as those mounted on the WorldView series.

8. Conclusion

In this paper, we have proposed the rigorous RPC warping model for the satellite MVS task for the first time. The advantage of the warping model is demonstrated both on the tests of the simulation data and the TLC SatMVS dataset. The experiments also show that the proposed learning based SatMVS framework performs better than the conventional method and the SOTA learning based methods when processing larger-capacity satellite images. In addition, we have presented the high-quality TLC SatMVS dataset. We believe that our work will promote the development of Earth surface reconstruction from MVS satellite images.

References

- [1] Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjarholm Dahl. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision*, 120(2):153–168, 2016. [2](#)
- [2] Marc Bosch, Kevin Foster, Gordon Christie, Sean Wang, Gregory D Hager, and Myron Brown. Semantic stereo for incidental satellite images. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1524–1532. IEEE, 2019. [2](#), [5](#), [8](#)
- [3] Marc Bosch, Zachary Kurtz, Shea Hagstrom, and Myron Brown. A multiple view stereo benchmark for satellite imagery. In *2016 IEEE Applied Imagery Pattern Recognition Workshop (AIPR)*, pages 1–9. IEEE, 2016. [2](#), [8](#)
- [4] Po-Heng Chen, Hsiao-Chien Yang, Kuan-Wen Chen, and Yong-Sheng Chen. Mvsnet++: Learning depth-based attention pyramid features for multi-view stereo. *IEEE Transactions on Image Processing*, 29:7261–7273, 2020. [1](#), [2](#), [3](#)
- [5] Shuo Cheng, Zexiang Xu, Shilin Zhu, Zhuwen Li, Li Erran Li, Ravi Ramamoorthi, and Hao Su. Deep stereo using adaptive thin volume representation with uncertainty awareness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2524–2534, 2020. [1](#), [2](#), [3](#), [4](#), [7](#)
- [6] Carlo De Franchis, Enric Meinhardt-Llopis, Julien Michel, Jean-Michel Morel, and Gabriele Facciolo. An automatic and modular stereo pipeline for pushbroom images. In *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, 2014. [2](#)
- [7] ASTER GDEM SRTM DTED. Aster global dem validation summary report. 2009. [5](#)
- [8] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361. IEEE, 2012. [2](#)
- [9] J. Grodecki and G. Dial. Block adjustment of high-resolution satellite images described by rational polynomials. *Photogrammetric Engineering and Remote Sensing*, 69:59–68, 2003. [3](#)
- [10] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping 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*, pages 2495–2504, 2020. [1](#), [2](#), [3](#), [4](#), [5](#), [7](#)
- [11] W. Hartmann, S. Galliani, M. Havlena, L. Van Gool, and K. Schindler. Learned multi-patch similarity. In *IEEE International Conference on Computer Vision (ICCV)*, pages 1595–1603, 2017. [2](#)
- [12] Heiko Hirschmüller. Stereo processing by semiglobal matching and mutual information. *IEEE Transactions on pattern analysis and machine intelligence*, 30(2):328–341, 2007. [2](#)
- [13] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2821–2830, 2018. [2](#)
- [14] inpho. Available: <https://geospatial.trimble.com/products-and-solutions/inpho>. [6](#)
- [15] Andy Jarvis, Hannes Isaak Reuter, Andrew Nelson, Edward Guevara, et al. Hole-filled srtm for the globe version 4. *available from the CGIAR-CSI SRTM 90m Database (<http://srtm.csi.cgiar.org>)*, 15:25–54, 2008. [5](#)
- [16] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. Surfacenet: An end-to-end 3d neural network for multiview stereopsis. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2307–2315, 2017. [2](#)
- [17] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, Peter Henry, Ryan Kennedy, Abraham Bachrach, and Adam Bry. End-to-end learning of geometry and context for deep stereo regression. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 66–75, 2017. [2](#)
- [18] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. [2](#)
- [19] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. [4](#)
- [20] Jin Liu and Shunping Ji. A novel recurrent encoder-decoder structure for large-scale multi-view stereo reconstruction from an open aerial dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6050–6059, 2020. [1](#), [2](#), [3](#), [4](#), [5](#), [7](#)
- [21] Keyang Luo, Tao Guan, Lili Ju, Haipeng Huang, and Yawei Luo. P-mvsnet: Learning patch-wise matching confidence aggregation for multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10452–10461, 2019. [1](#), [2](#), [3](#)
- [22] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4040–4048, 2016. [2](#)
- [23] Román Orús. A practical introduction to tensor networks: Matrix product states and projected entangled pair states. *Annals of Physics*, 349:117–158, 2014. [4](#)
- [24] Rongjun Qin. Rsp stereo processor (rsp)—a software package for digital surface model and orthophoto generation from satellite stereo imagery. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, 3:77, 2016. [2](#)
- [25] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. [4](#)
- [26] Daniel Scharstein and Richard Szeliski. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *International journal of computer vision*, 47(1):7–42, 2002. [2](#)

- [27] Thomas Schops, Johannes L Schonberger, Silvano Galliani, Torsten Sattler, Konrad Schindler, Marc Pollefeys, and Andreas Geiger. A multi-view stereo benchmark with high-resolution images and multi-camera videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3260–3269, 2017. [2](#)
- [28] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 519–528. IEEE, 2006. [2](#)
- [29] C Vincent Tao and Yong Hu. A comprehensive study of the rational function model for photogrammetric processing. *Photogrammetric engineering and remote sensing*, 67(12):1347–1358, 2001. [3](#)
- [30] Mi Wang, Fen Hu, and Jonathan Li. Epipolar resampling of linear pushbroom satellite imagery by a new epipolarity model. *ISPRS Journal of Photogrammetry and Remote Sensing*, 66(3):347–355, 2011. [2](#)
- [31] Jiayu Yang, Wei Mao, Jose M Alvarez, and Miaomiao Liu. Cost volume pyramid based depth inference for multi-view stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4877–4886, 2020. [1](#), [2](#), [3](#), [4](#)
- [32] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 767–783, 2018. [1](#), [2](#), [3](#), [5](#)
- [33] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5525–5534, 2019. [1](#), [2](#), [3](#), [5](#)
- [34] Kai Zhang, Noah Snavely, and Jin Sun. Leveraging vision reconstruction pipelines for satellite imagery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 2139–2148, 2019. [2](#), [5](#), [6](#), [7](#), [8](#)