

# Globally Consistent Alignment for Planar Mosaicking via Topology Analysis

Menghan Xia, Jian Yao\*, Renping Xie, Li Li

*Computer Vision and Remote Sensing (CVRS) Lab, School of Remote Sensing and Information Engineering, Wuhan University, Wuhan, Hubei, P.R. China*

Wei Zhang

*School of Control Science and Engineering, Shandong University, Jinan, Shandong, P.R. China*

---

## Abstract

In this paper, we propose a generic framework for globally consistent alignment of images captured from approximately planar scenes via topology analysis, capable of resisting the perspective distortion meanwhile preserving the local alignment accuracy. Firstly, to estimate the topological relations of images efficiently, we search for a main chain connecting all images over a fast built similarity table of image pairs (mainly for the unordered image sequence), along which the potential overlapping pairs are incrementally detected according to the gradually recovered geometric positions and orientations. Secondly, all the sequential images are organized as a spanning tree through applying a graph algorithm on the topological graph, so as to find the optimal reference image which minimizes the total number of error propagation. Thirdly, the global alignment under topology analysis is performed in the strategy that images are initially aligned by groups via the affine model, followed by the homography refinement under the anti-perspective constraint, which manages to keep the optimal balance between aligning precision and global consistency. Finally, experimental results on two challenging aerial image sets illustrate the superiority of the proposed approach.

*Keywords:* Topology Estimation, Reference Image, Graph Analysis, Global Consistency, Image Mosaicking

---

\*Corresponding author.

Email address: [jian.yao@whu.edu.cn](mailto:jian.yao@whu.edu.cn) (Jian Yao)

URL: <http://cvrs.whu.edu.cn/> (Jian Yao)

## **1. Introduction**

Nowadays, satellite and aerial remote sensing are common techniques to quickly capture images for territorial monitoring in both civil and military fields. Due to the limited observing range of a single image, image mosaicking is a necessary technique to stitch multiple images into a single wide-view seamless mosaic image. The first critical step in the mosaicking process is accurately aligning images into a common coordinate system, which directly influences the mosaicking quality [1, 2, 3]. As a strict aligning model, homography is often used to describe the relationship between two images of a 3D plane or two images captured from the same camera center [4, 5, 6, 7]. Recently, some mosaicking methods not limited to these two geometric conditions have been proposed to extend the range of applications [8, 9, 10]. Specially, in this paper, we focus on mosaicking images from an approximately planar scene. Challenged by both pseudo-plane and accumulation error, lots of related studies have been presented in the literature of the last decade. However, the performance of both accurate alignment and global consistency still remains to be further improved.

Generally, the image alignment approaches can be divided into two categories: area-based approaches [11, 12] and feature-based ones [13]. Because of the high computational cost, the area-based approaches are seldom used in the large scale mosaicking missions [14]. As for the planar mosaicking problem, such as aerial image mosaicking, the feature-based approaches are usually applied to recover the homography model between images [15, 16, 17] since the ground scene can be regarded as an approximate plane when it is observed from the aerial photographic camera. To improve the mosaicking result, many optimization algorithms have been proposed to achieve a global alignment. A typical global optimization method is “Bundle Adjustment” [18, 19], which aims at finding an optimal solution minimizing the total reprojection error [20]. To provide a good initial solution for global optimization, Xing et al. [21] proposed to first apply the Extended Kalman Filter [22] onto the local area, and then refine all the parameters globally. To avoid the non-linear optimization, Kekec et al. [23] employed the affine model to optimize the initial alignment recovered by the homography model in the global optimization. To prevent the image suffering down-scaling effect, Elibol et al. [24]

31 proposed to optimize point positions in the mosaicked frame and the alignment model in  
32 an alternate iteration scheme.

33 All those methods merely seek for an alignment with the least registration error, which  
34 can usually composite a satisfactory mosaic image from dozens of images. However, when  
35 sequential images are taken from a wide-range area, the global consistency of mosaicking  
36 images will be inaccessible for them, because in the case of pseudo-plane violating the  
37 strict geometric model, the least-registration-error principle is prone to inducing severe  
38 accumulation of perspective distortions. To avoid this problem, Caballero et al. [22]  
39 proposed to use the hierarchical models according to the alignment quality of images,  
40 where the model with the less degree of freedom (DoF) was used for images with the bigger  
41 parallax. The essence of this method is to make a trade-off between improving the aligning  
42 precision and resisting the perspective distortion. In fact, a more reasonable solution is to  
43 allow continuous transition between aligning models according the predefined constraint,  
44 which is detailed in our previous work [25].

45 As to the large-scale mosaicking problem, utilizing the topology among images is  
46 another effective way to improve the mosaicking result [26, 27]. To estimate the topology  
47 efficiently, Elibol et al. [28] used the low-cost tentative matching combined with the  
48 Minimum Spanning Tree (MST) solution to detect overlapping relations in an iterative  
49 scheme and decide when to update the topological estimation via the information-theory  
50 principles. As for the reference image selection, Richard et al. [29] stated that a reasonable  
51 choice is the most central image geometrically. This idea is obviously reasonable due to  
52 the fact that the central image usually has the shortest distance to all other images on  
53 average. To implement this idea, Choe et al. [30] applied a graph algorithm to select the  
54 reference image with the lowest cumulative registration error, but the registration error  
55 between each image pair have to be calculated in advance.

56 In this paper, we propose to obtain a visually satisfactory mosaic image with  
57 both accurate alignment and global consistency through two technical means: (1)  
58 utilizing the topology analysis to strengthen the registration constraints and reduce  
59 the error propagation; (2) adopting the alignment strategy of allowing continuous  
60 transition between different aligning models, to adaptively keep the optimal balance

61 between alignment accuracy and global consistency. Firstly, we initialize an approximate  
62 similarity matrix for image pairs efficiently, which is combined with the Minimum  
63 Spanning Tree (MST) to find the main chain for an unordered image sequence. Then,  
64 other potential overlapping relationships are detected incrementally with the gradually  
65 recovered geometric positions along the main chain. Because of the synchronism of overlap  
66 detection and image location, our topology estimation strategy is more efficient than the  
67 Elibol et al.’s method [28]. Secondly, all the sequential images are organized as a spanning  
68 tree through the classical Floyd-Warshall algorithm, so as to find the optimal reference  
69 image with the least cascading times. Finally, a globally consistent alignment strategy  
70 is applied, which combines the affine model with the homography model effectively. The  
71 initial alignment is recovered by the robust affine model by groups and the globally  
72 homography refinement is followed under the anti-perspective constraint. The proposed  
73 approach was tested by several experiments on two challenging aerial image datasets and  
74 the performances were comprehensively evaluated by comparing with the state-of-the-art  
75 algorithm and a famous commercial software.

76 The remainder of this paper is organized as follows. The proposed framework is  
77 detailed in Section 2, which is comprised of the topology estimation, the selection of  
78 reference image, and the global alignment. Experimental results are provided in Section 3,  
79 followed by the conclusion and future work presented in Section 4.

## 80 2. Our Approach

81 We propose a generic framework for globally consistent alignment of images captured  
82 from an approximately planar scene as shown in Figure 1, which includes three modules:  
83 topology estimation, reference image selection, and global alignment. First, the sequential  
84 images are inputed for topology estimation, through which the obtained topological graph  
85 and matching results are utilized to search the optimal reference image and to provide  
86 feature correspondences for the global alignment respectively. Finally, according to the  
87 reorganized hierarchy, all the images are aligned through a specially designed double-  
88 model optimization strategy. Due to the versatile topology estimation, the proposed  
89 framework is suitable for both time-consecutive image sequence and unordered one.

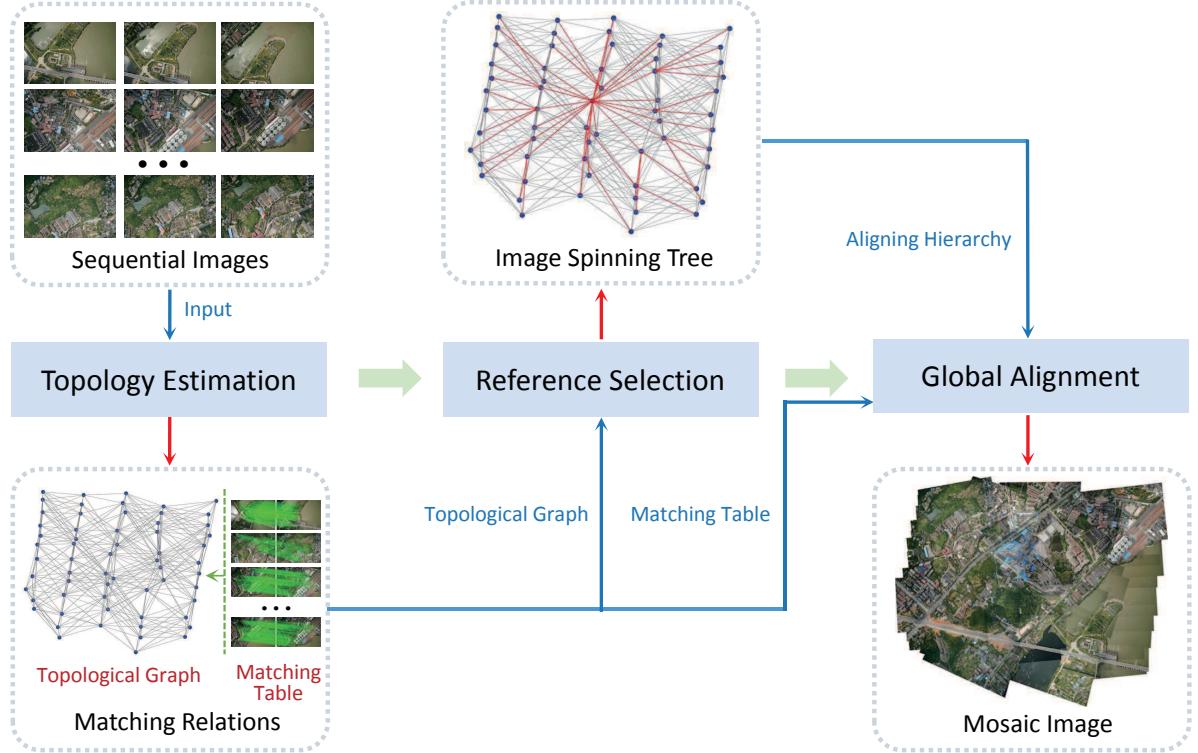


Figure 1: The flowchart of our proposed framework for globally consistent alignment of images. The blue and red thin arrows denote the input and output of each processing module, respectively, and the wide green arrows indicate the execution sequence.

For the description convenience in the following, the frequently used notations in this paper are summarized below :

- $\mathbf{I}_i$  - the  $i$ -th image in the sequential images.
- $\mathbf{A}_i$  - the  $3 \times 3$  affine transformation matrix relating  $\mathbf{I}_i$  to the reference frame.
- $\mathbf{H}_i$  - the  $3 \times 3$  homography transformation matrix relating  $\mathbf{I}_i$  to the reference frame.
- $\mathbf{x} = [x, y, w]^\top$  - the homogeneous coordinate of a 2D image feature point.
- $\mathbf{x}_{i,j}^k$  - the 2D coordinate of the  $k$ -th matched feature in  $\mathbf{I}_i$  corresponding to the  $k$ -th matched feature  $\mathbf{x}_{j,i}^k$  in  $\mathbf{I}_j$ .
- $M_{i,j}$  - the total number of feature matches between  $\mathbf{I}_i$  and  $\mathbf{I}_j$ .
- $\varpi(\mathbf{x}) = [x/w, y/w]^\top$  - the function transforming the homogeneous coordinate of a 2D point into the non-homogeneous coordinate.

101    *2.1. Fast Topology Estimation*

102    The image topology of the surveyed area is usually represented by a graph where  
103    each image is depicted as a node and the overlapping relationship between image pair is  
104    denoted by an edge or a link. Topology estimation means to find the existing overlapping  
105    relationships among all images. In this section, we try to find all the potential overlapping  
106    image pairs by utilizing the gradually recovered geometric positions of images in the time-  
107    consecutive order on the mosaicking plane, instead of blindly doing matching attempts.  
108    As for an unordered image sequence, finding a main chain connecting all images can  
109    make the problem the same as that of the time-consecutive image sequence. Therefore,  
110    an efficient strategy can be proposed to find the complete topology with the minimum  
111    image matching attempts.

112    *2.1.1. Finding Main Chain with Most Reliability*

113    For a sequence of  $n$  images, the main chain consisting of  $(n - 1)$  edges connects all the  
114    nodes/images in the graph. More strictly, it is defined as a spanning tree of an undirected  
115    graph in graphic theory [31, 32]. Usually, there is no need to find a main chain for a time-  
116    consecutive image sequence since their time-consecutive links have implied a main chain  
117    already. Thus, this step is mainly set for an unordered image set in topology estimation.

118    Given an unordered image set, we have to measure the similarities between image pairs  
119    in advance for finding a main chain. Here, the initial similarity information is intended  
120    to be computed in an approximate and efficient way. For each image, we first extract  
121    SURF [33] features and only select a particular feature subset to represent this image.  
122    Then, the similarity between image pair is defined as the number of candidate point  
123    matches whose descriptor vector distances are less than some given distance. Specially,  
124    to increase the corresponding probability, the feature subset of every image consists of  
125    the features extracted on the same scale layer defined in the SURF detector, instead of  
126    sampling randomly. In our approach, the features from the second scale layer of the total  
127    four octaves are selected as the subset representing each image, because features from  
128    this layer hold a stable ratio of  $22 \pm 3\%$  (an appropriate sampling ratio) almost for all  
129    kinds of images. The computational cost of obtaining the initial similarity information is  
130    comparatively low, since it mainly involves computing the distances between descriptor

vectors of feature subsets. Over the exhaustive comparison, all the similarity values between image pairs are organized in the form of a matrix  $\mathbf{S}$ , where  $\mathbf{S}(i, j)$  represents the similarity between images  $\mathbf{I}_i$  and  $\mathbf{I}_j$ . The value of  $\mathbf{S}(i, j)$  from small to large means an increasing similarity between images  $\mathbf{I}_i$  and  $\mathbf{I}_j$ , which can be regarded as the probability of images  $\mathbf{I}_i$  and  $\mathbf{I}_j$  sharing an overlap.

Based on the similarity matrix, the reciprocals of those non-zero similarity values are set as the weights for the edges of the graph, i.e.,  $W(i, j) = \frac{1}{\mathbf{S}(i, j)}$ . Considering the similarity is not reliable, we try to find a valid main chain through an iterative scheme between selecting the main chain candidate and verifying its connectivity:

**Maximum Reliability:** Given such a weighted graph, we try to select a linkage path that connects all the nodes with the highest total reliability, i.e., the lowest sum of weights. This is realized by finding the Minimum Spanning Tree (MST) of the current weighted graph. The MST is a spanning tree whose edges have the minimum total weight in all the spanning trees of the graph. So, the edges of MST represent the most possible overlapping relationships between image pairs.

**Connectivity Verification:** To verify the real connectivity depicted by edges of the MST, a more reliable feature matching algorithm is employed to check the overlapping relationships between these image pairs. In such matching, all the SURF features are used and both epipolar constraint and appropriately homography constraint are applied to remove outliers. If all the matching attempts succeed, the MST is a valid main chain and the iteration terminates. Reversely, when there exists any failed matching attempt, we modify the weights of the graph where the weights of successfully matched pairs are set as zero while the weights of matching-failed image pairs are set as an infinite value, then it turns to the next iteration.

Specifically, when it fails to find the MST or the iteration reaches a given threshold, it means no connected graph exists, and the procedure quits. To illustrate the property of main chain, a subset of the first dataset described in Section 3 was selected to demonstrate the results of topology estimation, which was tested in the time-consecutive mode and in the unordered mode, respectively, as shown in Figure 2. Apart from the difference of the main chains, the topologies estimated in the two different modes are almost the same,

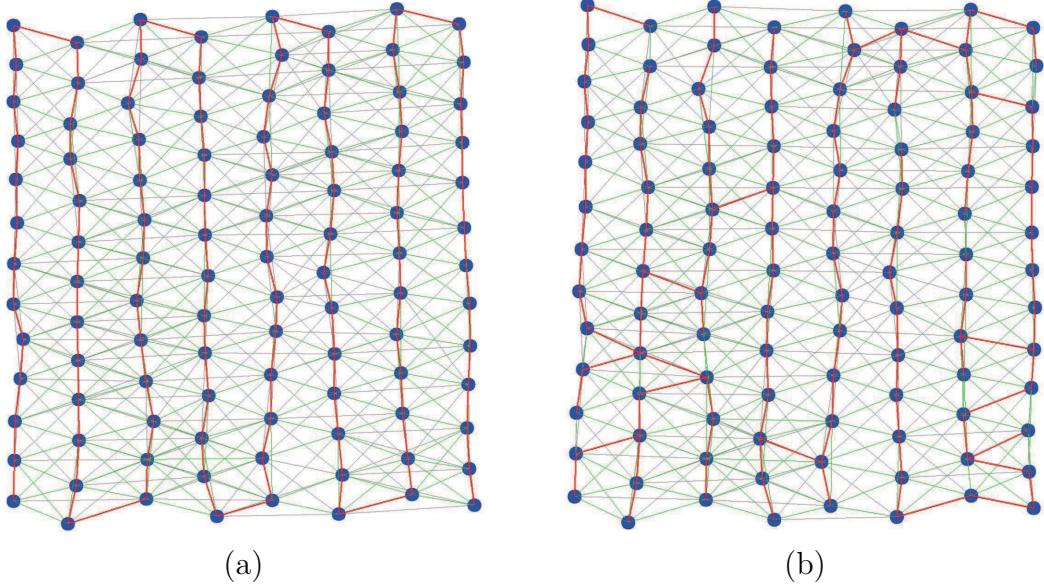


Figure 2: The estimated topologies of an image sequence (104 images) in the time-consecutive mode and the unordered one, respectively: (a) the topology estimated in the time-consecutive mode, where the red edges represent the prior main chain in the time-consecutive order; (b) the topology estimated in the unordered mode, where the red edges represent the main chain linked by the proposed iterative scheme. Besides, the green edge indicate more than 100 matches between an image pair, while the gray edge indicate the number of the matches are less than 100.

161 which are compared quantitatively in Table 1.

### 162 2.1.2. Detecting Potential Overlapping Pairs

163 According to overlapping relationships indicated by the main chain, we can recover the  
 164 comparative geometric positions of sequential images by projecting them into a common  
 165 coordinate system. Then, based on the geometric information, the potential overlapping  
 166 pairs can be detected easily. In our approach, the two operations are conducted in a  
 167 synchronously collaborative way, instead of an independently serial way.

168 Firstly, we temporarily select a reference image as the mosaicking plane through  
 169 applying the algorithm detailed in Section 2.2 on the main chain. To recover the  
 170 comparatively geometric positions, we employ the affine model to align images, which  
 171 is robust in locating the centroids of images. Compared with the affine model, the  
 172 homography model is prone to suffering from the perspective distortion and the 2D rigid  
 173 model tends to cause a bending trajectory because of the error accumulation, which was  
 174 testified in Section 3.2. To improve the reliability of the image locations, the images on the

---

**Algorithm 1** Detecting potential overlapping pairs

---

**Input:** The image set  $\mathcal{I} = \{\mathbf{I}_i\}_{i=1}^n$  arranged in the order of breadth-first searching the main chain spanning tree (with the reference image as the root node).

**Output:** The set of overlapping image pairs  $\mathcal{P} = \{\mathbf{P}_{ij}\}_{i \neq j}$ .

- 1: Initialize the located image set  $\widehat{\mathcal{I}} = \{\mathbf{I}_1\}$ .
- 2: **for** each image  $\mathbf{I}_i \in \mathcal{I} \setminus \{\mathbf{I}_1\}$  **do**
- 3:     Align  $\mathbf{I}_i$  with its direct reference image  $\mathbf{I}_{\rho(i)}$ .
- 4:     Initialize the overlapping pairs set  $\mathcal{P}_i = \{\mathbf{P}_{i\rho(i)}\}$ .
- 5:     **for** each image  $\mathbf{I}_j \in \widehat{\mathcal{I}} \setminus \{\mathbf{I}_{\rho(i)}\}$  **do**
- 6:         yes/no  $\leftarrow$  Detect the overlap between  $\mathbf{I}_i$  and  $\mathbf{I}_j$ .
- 7:         **if** yes **then**
- 8:              $\mathcal{P}_i = \mathcal{P}_i \cup \{\mathbf{P}_{ij}\}$ .
- 9:         **end if**
- 10:       **end for**
- 11:     Realign  $\mathbf{I}_i$  with its neighborhood image set  $\mathcal{P}_i$ .
- 12:      $\mathcal{P} = \mathcal{P} \cup \mathcal{P}_i$
- 13:      $\widehat{\mathcal{I}} = \widehat{\mathcal{I}} \cup \{\mathbf{I}_i\}$
- 14: **end for**
- 15: **return**  $\mathcal{P}$

---

175 main chain will be aligned starting from the reference image one by one. As the images are  
176 located gradually, the potential overlapping relationships around the newly located image  
177 are detected, which then is used for optimizing the position of this newly aligned image.  
178 This strategy benefits in improving the accuracy of the recovered geometric positions,  
179 because the simultaneously detected overlapping pairs can provide extra constraints for  
180 alignment. Given a newly aligned image  $\mathbf{I}_i$ , we check whether it shares an overlap with  
181 all the previously aligned images  $\widehat{\mathcal{I}} = \{\mathbf{I}_j\}_{j=1}^m$ . Specially, the overlap detection between  
182  $\mathbf{I}_i$  and  $\mathbf{I}_j$  is performed by calculating the distance between their centroids as follows:

$$\delta_{ij} = \frac{\max(0, |c_i - c_j| - |d_i - d_j|/2)}{\min(d_i, d_j)}, \quad (1)$$

183 where  $c_i$ ,  $c_j$ ,  $d_i$  and  $d_j$  are the image centroids and the diameters of the minimum boundary  
184 circles of the projection onto the mosaicking plane of  $\mathbf{I}_i$  and  $\mathbf{I}_j$ , respectively. If  $\delta_{ij} > 1$ ,

Table 1: Comparisons of our topology estimation running in both the time-consecutive mode (a) and the unordered mode (b) (with All-against-all as the ground truth).

Strategy	Successful Attempts	Total Attempts	% of Recall	% of Computation on Feature Matching
Proposed Approach (a)	606	896	94.71	99.42
Proposed Approach (b)	595	905	92.10	84.14
All-against-all	646	5356	100.00	100.00

185 there is no overlap. Otherwise, there may exist an overlap between  $\mathbf{I}_i$  and  $\mathbf{I}_j$ , and we  
 186 attempt to match them for verification. Of course, if the matching between  $\mathbf{I}_i$  and  $\mathbf{I}_j$  has  
 187 been done in the stage of finding the main chain, it is no need to do the matching attempt  
 188 again. The procedure of our topology estimation approach is described in Algorithm 1.  
 189 When all the overlapping pairs are obtained, we redefine the similarity matrix as the final  
 190 topological representation. The original similarity matrix is reset as a zero matrix firstly,  
 191 and the value of  $\mathbf{S}(i, j)$  is replaced with the number of matched points only if  $\mathbf{I}_i$  and  $\mathbf{I}_j$   
 192 have been matched successfully.

193 It should be noted that the major computation cost of the topology estimation is on  
 194 feature matching between images, as listed in the fourth column of Table 1. The image  
 195 alignment and potential overlapping detection have a relatively low computation cost,  
 196 since there is no global optimization or iterative detection involved. Besides, as for an  
 197 unordered image set, the initialization of the similarity matrix occupies the majority of  
 198 the rest computation.

### 199 *2.2. Optimal Reference Image Selection*

200 As we known, the images alignment is realized through warping each image onto  
 201 the mosaicking plane which is usually one of the input images (named as the reference  
 202 image). Projecting an image without direct overlap with the reference image to the  
 203 mosaicking plane involves cascading a series of relative transformation models of other  
 204 intermediate images. Obviously, less intermediate images used for cascading makes less  
 205 error accumulation. In fact, there might exist more than one path of the same cascading  
 206 number from an image to another. Considering each cascading induces a different error,

we manage to select the path with the least accumulation error. In terms of this, the optimal reference image should give the lowest sum of accumulation errors from all the other images to the reference image plane. To address this problem, we construct an undirected weighted graph based on the estimated topology in Section 2.1. According to the final similarity matrix, image pairs with non-zero values of similarities are linked with edges. As to the weight (or cost) of an edge, there are two kinds of settings in the existing literature: the reciprocal of the number of matched features [28] and the registration error between the image pair [30]. The former is intuitive and efficient while the latter depicts the error directly at the cost of calculating the registration error between all available image pairs in advance. Considering the association between the number of matched features and the registration error, we creatively set the edge weight as follows:

$$w_{ij} = \begin{cases} \inf, & \text{if } M_{i,j} = 0, \\ \frac{1}{\log(M_{i,j} + \varepsilon)}, & \text{if } M_{i,j} > 0, \end{cases} \quad (2)$$

where  $M_{i,j}$  denotes the total number of matches between  $\mathbf{I}_i$  and  $\mathbf{I}_j$ , and  $\varepsilon$  is a constant for regularization ( $\varepsilon = 50$  by default). This weight setting equation, which describes the contribution of matched features to the registration accuracy, compromises between efficiency and effectiveness.

Based on the weighted graph, the optimal reference image selection problem is formulated as finding a node with the least total weight of the shortest paths to all the other nodes, which can be solved by the Floyd Warshalls all-pairs shortest path algorithm [34, 35]. This algorithm is more efficient than running  $n$  times of single source shortest path algorithm, because the dynamic programming strategy is applied in it with the computation complexity of  $O(n^3)$ . With this algorithm, the shortest paths between any two nodes can be obtained. For a sequence of  $n$  images, we build a  $n \times n$  size symmetric cost matrix  $\mathbf{W}$  where  $\mathbf{W}(i, j)$  records the cost of the shortest path between  $\mathbf{I}_i$  to  $\mathbf{I}_j$ . Thus, the  $i$ -th row or column of matrix  $\mathbf{W}$  indicates the cost of every shortest path from other images to  $\mathbf{I}_i$ , and the column with the minimum cumulative cost is selected as the reference image. To demonstrate the procedure, the cost matrix  $\mathbf{W}$  of a sequence of 104 images is visualized in Figure 3. The 45-th column with the minimum total cost is marked with a red arrow in the bottom indices. As a comparison, the conventional

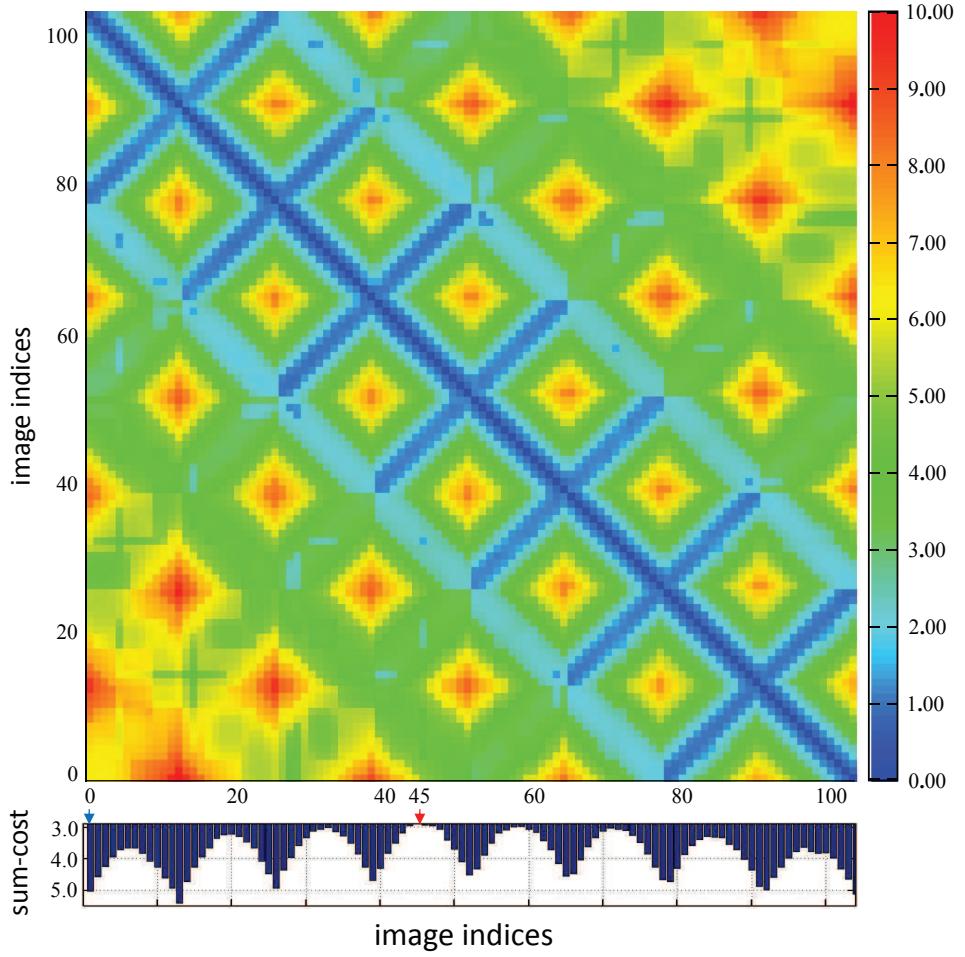


Figure 3: The cost matrix of all-pairs shortest path calculated from a sequence of 104 images. Below is attached by the bar chart depicting the mean cumulative cost of each column of the cost matrix. A red arrow in the bottom indices labels the 45-th column with the minimum mean cumulative cost of 3.04. For comparison, the first column labeled with a blue arrow has a much higher mean cumulative cost of 5.23.

strategy of naively selecting the first image as the reference image is marked with a blue arrow. Considering the amount of images, the gap between the mean cumulative costs of the two strategies can make a big difference to the mosaicking result.

Actually, each row in  $\mathbf{W}$ , e.g. the  $i$ -th row, corresponds to a spanning tree with the node  $i$  as the root one, which describes the hierarchical relationships of all nodes. With the selected reference image, the spanning tree of the image sequence used in Figure 3, is displayed in Figure 4. The spanning tree indicates the direct reference image (i.e., the parent node) of each image, which determines the aligning order of images in the

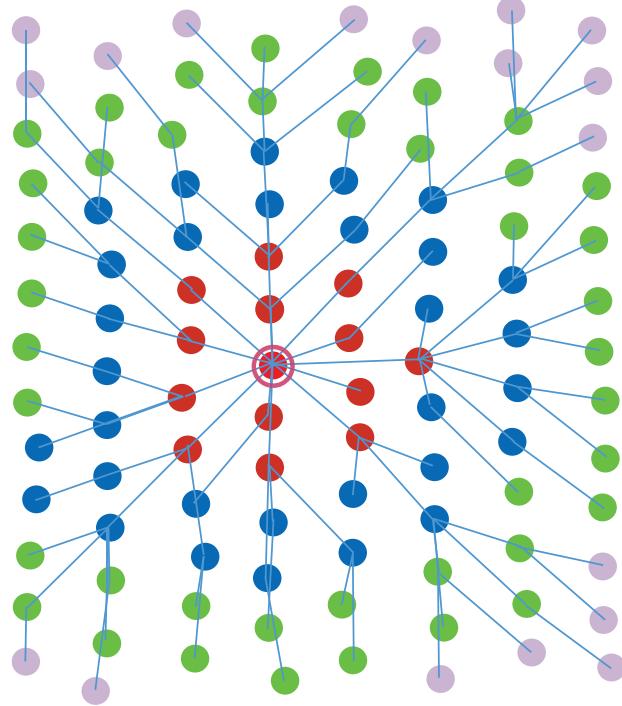


Figure 4: The spanning tree of the graph with the optimal reference image as the root node (marked with red ring). Nodes in different levels of the tree are marked with different colors, and the blue lines imply the shortest path from every image to the reference image.

<sup>243</sup> following global alignment.

<sup>244</sup> *2.3. Globally Consistent Alignment*

<sup>245</sup> In general, the local aligning accuracy and the global consistency are two basic factors  
<sup>246</sup> determining the quality of mosaicking result. Under a strict transformation model, these  
<sup>247</sup> two factors can be guaranteed in a coherent way, where the higher aligning precision  
<sup>248</sup> contributes on the better global consistency. However, in most applications, the observing  
<sup>249</sup> scenes of pesudo-planes make the frequently-used homographic transformation just an  
<sup>250</sup> approximate aligning model. In this case, the model of higher degrees of freedom (DoFs)  
<sup>251</sup> usually makes more accurate alignment but suffers more severe perspective distortion  
<sup>252</sup> meanwhile, and vice versa. Therefore, we have to deal with these two factors in a trade-  
<sup>253</sup> off way. To keep the optimal balance between them, the model with a relatively low DoFs  
<sup>254</sup> is employed for initial alignment which is robust to the perspective distortion, it then is  
<sup>255</sup> refined with a higher DoFs to improve the aligning precision under the anti-perspective  
<sup>256</sup> constraint.

257    2.3.1. Robust Alignment by Affine Model

258    The affine model, compromising between the 2D rigid transformation and the  
 259    homographic transformation, is used to make a robustly initial alignment. On the  
 260    one hand, the approximately coplanar constraint of images is partly implied in the six-  
 261    parameter affine model which can suppress the perspective distortion in some extent, on  
 262    the other hand, the affine transformation is able to provide a qualified initial solution for  
 263    the following homography refinement.

264    According to the spanning tree mentioned in Section 2.2, the sequential images  
 265    are aligned group by group in the order of breadth-first search, which reduces the  
 266    accumulation error of alignment compared to the way one by one. When aligning a  
 267    new group of images to the reference frame, the overlapping relations between all the  
 268    previously aligned images and the newly added ones, and the overlapping relations inside  
 269    this group will be jointly used in the cost function. Let  $\mathcal{I} = \{\mathbf{I}_i\}_{i=1}^s$  be the set of previously  
 270    aligned images. The affine model parameters  $\mathcal{A} = \{\mathbf{A}_i\}_{i=s+1}^{s+m}$  of the newly added image  
 271    group  $\mathcal{G} = \{\mathbf{I}_i\}_{i=s+1}^{s+m}$  will be solved by minimizing the cost function as below:

$$E(\mathcal{A}) = E_1(\mathcal{A}|\mathcal{I}, \mathcal{G}) + E_2(\mathcal{A}|\mathcal{G}), \quad (3)$$

272    where the first energy term  $E_1(\mathcal{A}|\mathcal{I}, \mathcal{G})$  corresponds to the overlapping relations between  
 273     $\mathcal{I}$  and  $\mathcal{G}$  as:

$$E_1(\mathcal{A}|\mathcal{I}, \mathcal{G}) = \sum_{\mathbf{I}_i \in \mathcal{I}, \mathbf{I}_j \in \mathcal{G}} \sum_{k=1}^{M_{i,j}} \|\varpi(\mathbf{A}_i \mathbf{x}_{i,j}^k) - \varpi(\mathbf{A}_j) \mathbf{x}_{j,i}^k\|^2, \quad (4)$$

274    and the second energy term  $E_2(\mathcal{A}|\mathcal{G})$  corresponds to the intra-group overlapping relations  
 275    in  $\mathcal{G}$  as:

$$E_2(\mathcal{A}|\mathcal{G}) = \sum_{\mathbf{I}_i, \mathbf{I}_j \in \mathcal{G}} \sum_{k=1}^{M_{i,j}} \|\varpi(\mathbf{A}_i \mathbf{x}_{i,j}^k) - \varpi(\mathbf{A}_j) \mathbf{x}_{j,i}^k\|^2, \quad (5)$$

276    where the meanings of the notations  $\varpi(\cdot)$ ,  $\mathbf{A}_i$ ,  $M_{i,j}$  and  $\mathbf{x}_{i,j}^k$  are given in the beginning of  
 277    Section 2.

278    As the linear equations, Eq. (3) can be solved easily by the Singular Value  
 279    Decomposition (SVD) method. To increase the numerical solution stability, we normalize  
 280    the coordinates of matched points [36] to build the coefficient matrix. In addition, the  
 281    robust estimator MLESAC [37] is used to exclude outliers for affine estimation because  
 282    it is beneficial for the image mosaicking of quasi-planar scenes.

283    2.3.2. Model Refinement under Anti-Perspective Constraint

284    The affine models recovered by groups are mainly used to make the robust initial  
 285 alignment, which well preserves the mosaicking result from the perspective distortion.  
 286 However, due to the fact that the DoF of the model is limited and no global optimization  
 287 is performed, it is necessary and probable to improve the aligning precision further. To  
 288 increase the aligning accuracy to the extent not inducing the perspective distortion, the  
 289 energy function should allow the transition from the affine model to the homography  
 290 model under some reasonable constraint. In fact, such constraint has been implied in  
 291 the affine model which has the anti-perspective property. So, the deviation between the  
 292 optimal homography transformation and the initially estimated affine transformation is  
 293 set as a regularization term in the optimization function.

294    As the images are aligned by groups, the affine models of all the images  $\mathcal{I} = \{\mathbf{I}_i\}_{i=1}^n$   
 295 can be obtained, denoted as  $\mathcal{A} = \{\mathbf{A}_i\}_{i=1}^n$ , which are used as the initial parameters for the  
 296 homography models in the global optimization. The homography models  $\mathcal{H} = \{\mathbf{H}_i\}_{i=1}^n$   
 297 with respective to the reference plane will be optimized in the energy function composed  
 298 of two mutually contrary terms. The data term targeting to minimize the sum of squares  
 299 of the feature registration errors between images is defined as:

$$E_d(\mathcal{H}) = \sum_{\mathbf{I}_p, \mathbf{I}_q \in \mathcal{I}} \sum_{k=1}^{M_{p,q}} \|\varpi(\mathbf{H}_p \mathbf{x}_{p,q}^k) - \varpi(\mathbf{H}_q \mathbf{x}_{q,p}^k)\|^2, \quad (6)$$

300 where all the aligning models have more free parameters to adjust the positions of points  
 301 on the mosaicking plane, which is bound to increase the whole alignment precision.  
 302 Besides, the residual error is prone to distributing evenly under an uniform energy  
 303 framework.

304    Another optimization objective is to keep the global consistency by suppressing the  
 305 accumulation of the perspective distortions which may emerge in the transition from the  
 306 affine model to the homography model. The regularization term, derived from the idea  
 307 that the optimal homography transformation should be close to the initially estimated  
 308 affine transformation, is expressed as the displacements of the warped features from their  
 309 initial positions as:

$$\begin{aligned} E_r(\mathcal{H}) = & \sum_{\mathbf{I}_p, \mathbf{I}_q \in \mathcal{I}} \sum_{k=1}^{M_{p,q}} (\|\varpi(\mathbf{H}_p \mathbf{x}_{p,q}^k) - \mathbf{A}_p \mathbf{x}_{p,q}^k\|^2 \\ & + \|\varpi(\mathbf{H}_q \mathbf{x}_{q,p}^k) - \mathbf{A}_q \mathbf{x}_{q,p}^k\|^2). \end{aligned} \quad (7)$$

310 As depicted in Eq. (7), the regularization term is also denoted by the distances of  
311 image feature points as the data term does, which saves the troublesome normalized  
312 problem between different kinds of energy terms. So far, the energy terms defined in  
313 Eq. (6) and Eq. (7) can be linearly combined to define the final energy function as:

$$E(\mathcal{H}) = E_d(\mathcal{H}) + \lambda E_r(\mathcal{H}), \quad (8)$$

314 where  $\lambda$  is the weight coefficient used for balancing the two terms  $E_d$  and  $E_r$ , which should  
315 be set to an appropriate small value since the constraint isn't a strict one. Theoretically,  
316 a bigger value of  $\lambda$  strengthens the global consistency while decreases the accuracy of  
317 the local alignment. For instance, for images of large-depth-difference ground, a slightly  
318 bigger  $\lambda$  is needed to resist perspective distortion. In our experiments, we set its value  
319 from 0.01 to 0.05 according to data traits. As a typical non-linear least squares problem,  
320 Eq. (8) can be solved by the Levenberg-Marquardt (LM) algorithm. However, considering  
321 the specialty of this problem, we employ the sparse LM algorithm [38] to save memory  
322 and to speed up the computation, which is stated detailedly in Appendix A.

### 323 3. Experimental Results

324 In this section, two sets of representative aerial images acquired by different flight  
325 platforms and over different landforms, respectively, were used as the experimental  
326 dataset. The first dataset, consisting of 744 images from 24 sequentially ordered strips,  
327 was captured at a flight height of about 780 meters over an urban area. The original  
328 images with a forward overlapping rate of about 60%, are down-sampled to the size of  
329  $1000 \times 642$  in our experiments. The second dataset, consisting of 130 images with the  
330 down-sampling size of  $800 \times 533$ , was captured by an unmanned aerial vehicle (UAV)  
331 with a forward overlapping rate of about 70%, which observes a suburb area containing  
332 mountains.

333 Due to the limit of pages, more experimental results and analysis are presented at  
334 <http://cvrs.whu.edu.cn/projects/PlanarMosaicking/>, where the dataset and the  
335 source code are publicly available for download.

336    *3.1. Evaluation on Topology Estimation*

337    Our topology estimation approach was compared with the classic all-against-all  
338    strategy and the state-of-the-art algorithm implemented according to [28] (we name it  
339    as Fast-Topology hereafter). The comparisons were made on the estimated topologies of  
340    the aforementioned datasets. To test our approach diversely, the aerial image sequence  
341    and the UAV image sequence were respectively processed in two different modes for  
342    topology estimation: the time-consecutive mode and the unordered mode. As a robust  
343    but exhaustive strategy, matching all-against-all can give the topology estimation result  
344    which can be regarded as the ground truth. Moreover, the successfully matched image  
345    pairs and the total matching attempts were combined to evaluate the topology estimation  
346    algorithm as quantitative metrics in accuracy and efficiency.

347    The topology estimation results of the two datasets are summarized in Table 2 and  
348    Table 3, respectively. As the tables show, both our approach and Fast-Topology [28]  
349    almost recovered the complete topology as the all-against-all strategy does, but with  
350    much less matching attempts. Although there are some omissions with respect to  
351    all-against-all, the major overlapping relations have been detected successfully in our  
352    approach, which can be observed in the topological graph depicted in Figure 5(a) and  
353    Figure 6(a). This implies that most of the undetected overlapping pairs probably share  
354    very small overlapping areas, and usually make little contribution to the mosaicking  
355    results. Compared to Fast-Topology [28], our approach has roughly the same recall rates  
356    but less total matching attempts, which benefits from two key strategies used in the  
357    potential overlapping pairs detection. The one is selecting a suitable temporary reference  
358    image by applying the strategy detailed in Section 2.2, instead of selecting the first image  
359    simply as in Fast-Topology. The other is that the position of the newly added image  
360    is simultaneously adapted along with the potential overlapping relations being detected,  
361    which improves the alignment accuracy and so does the efficiency. However, in Fast-  
362    Topology, detecting the potential overlapping pairs and adapting alignment of images  
363    with the detecting results are divided into two independent steps. Thus, it inevitably  
364    introduces many unnecessary matching attempts because of the inaccurate alignment in  
365    the first few iterations, though it can find most of the existing overlapping relations after

Table 2: Comparisons of the topology estimation obtained by different approaches on the first dataset (with All-against-all as the ground truth).

Strategy	Successful Attempts	Total Attempts	% of Recall	% of Attempts As to All-against-all
Our Approach	5197	7771	97.83	2.81
Fast-Topology [28]	5229	9601	98.43	3.47
All-against-all	5312	276396	100.00	100.00

Table 3: Comparisons of the topology estimation obtained by different approaches on the second dataset (with All-against-all as the ground truth).

Strategy	Successful Attempts	Total Attempts	% of Recall	% of Attempts As to All-against-all
Our Approach	781	934	95.36	11.14
Fast-Topology [28]	793	1336	96.83	15.93
All-against-all	819	8385	100.00	100.00

366 several iterations.

367 As mentioned in Section 2.2, the estimated topology is used to search for the optimal  
368 reference image, by the way of which the images are organized as a spanning tree implying  
369 the aligning order for the global alignment. Here, the spanning trees with the reference  
370 image as the root node, are expressed by a group of red edges of the topological graph in  
371 Figure 5(b) and Figure 6(b), corresponding to the first and second datasets, respectively.  
372 It's easy to find that the selected reference images can always locate in the central part  
373 geometrically, no matter of the square-shaped aerial data or the strip-shaped UAV data.

374 *3.2. Evaluation on Initial Model Selection*

375 In the period of recovering initial alignment described in Section 2.3.1, the selection  
376 of the transformation model among *rigid*, *affine* and *homography* models can make  
377 differences to the final mosaicking result. To amplify the influence of error factors, we  
378 specially selected a strip-shaped aerial image subset and a block UAV image subset from  
379 the first dataset and the second one, respectively, and the image on the end was set as  
380 the reference image. The comparative analyses were made on both alignment precision

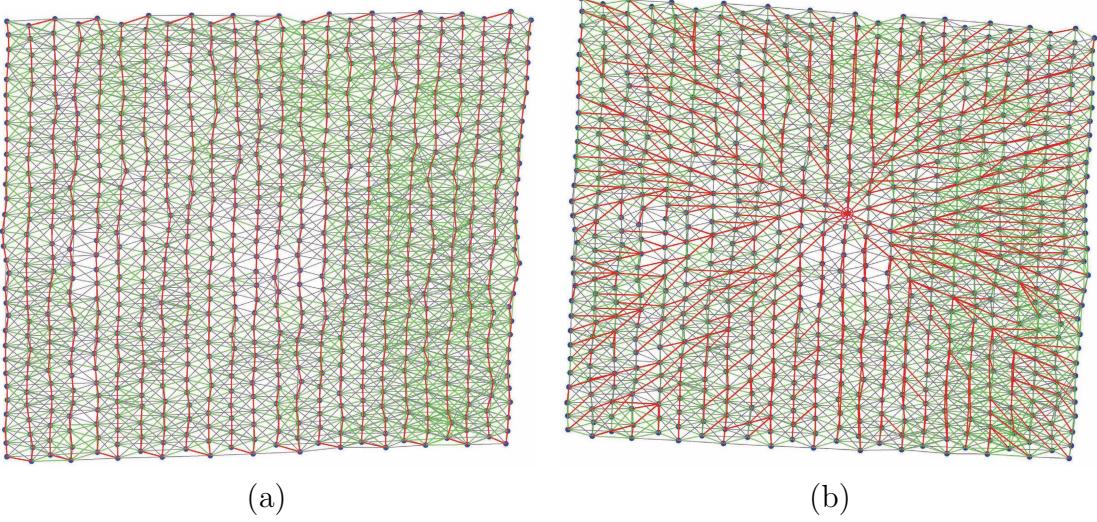


Figure 5: The estimated topology of the first dataset (744 images) highlighted for different aims: (a) the estimated topology with the prior main chain marked with red edges; (b) the spanning tree generated by searching for the optimal reference image, marked with red edges on the estimated topological graph. Different from (a), the geometric positions in (b) are recovered by the final global alignment. The edge in green and gray indicate the matches between the image pair more and less than 100 respectively.

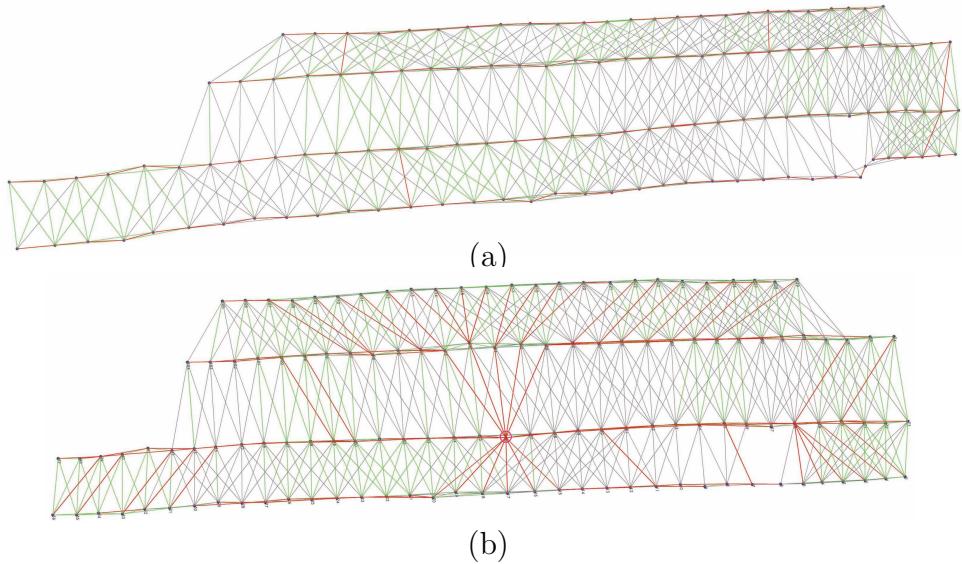


Figure 6: The estimated topology of the second dataset (130 images) highlighted for different aims: (a) the estimated topology with the searched main chain marked with red edges; (b) the spanning tree generated by searching for the optimal reference image, marked with red edges on the estimated topological graph. Different from (a), the geometric positions in (b) are recovered by the final global alignment. The edge in green and gray indicate the matches between the image pair more and less than 100 respectively.

Table 4: The Root-Mean-Square (RMS) errors through selecting different transformation models for initial alignment in the proposed approach (GR: Global Refinement; Unit: pixel).

Models	Strip Aerial Images			Block UAV Images		
	#Matches	RMS	RMS (GR)	#Matches	RMS	RMS (GR)
Rigid	131279	3.142	1.247	48783	5.112	1.985
Affine	131279	2.825	1.117	48783	4.421	1.743
Homography	131279	2.459	0.808	48783	3.605	1.485

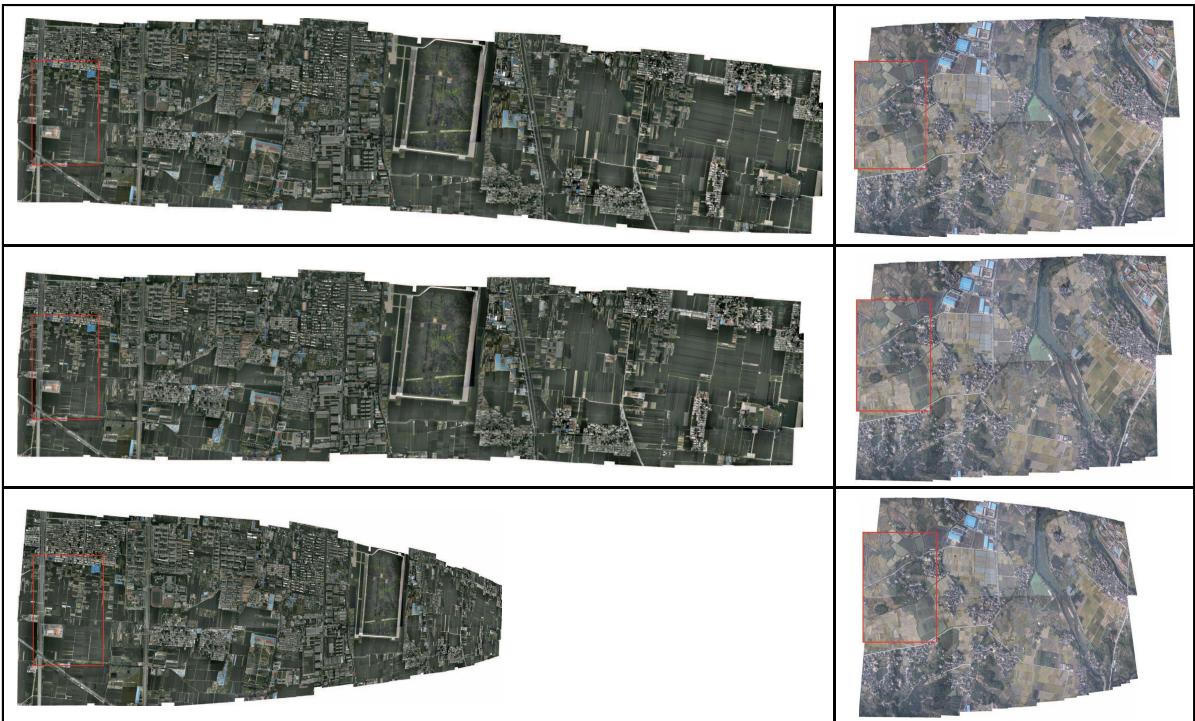


Figure 7: The thumbnails of the mosaicking results on the aerial images (Left) and the UAV images (Right) where the rigid model in the first row, the affine model in the second row, and the homography model in the last row are chosen for initial alignment, respectively. Notice that the reference image of each mosaic is marked with a red rectangular box.

and global consistency, where the numerical results are shown in Table 4 while the global consistency can be judged via the visual results shown in Figure 7.

As for the strip aerial images, the homography model as the initial model has the best alignment precision, but suffers severe an accumulation of the perspective distortions because it has the highest DoF for alignment. However, the mosaicking result of rigid

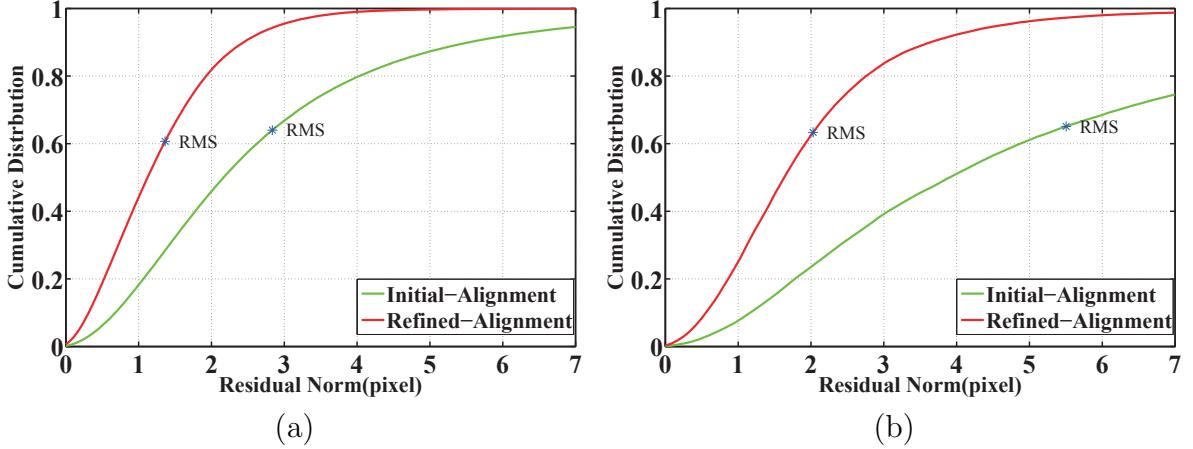


Figure 8: Cumulative probability distributions of the residual error norms with and without the global refinement performed in our approach: (a) the error analysis for the first dataset; (b) the error analysis for the second dataset. The green curve depicts the precision of initial alignment, while the red curve depicts the precision of global alignment. The blue marks on curves indicate the RMS errors.

transformation shows a bending tendency with the lowest accuracy although it doesn't induce any perspective distortion. This is because the rigid model of 3 free parameters just allows the image translation and rotation, which is not enough to describe the truly geometric relations and easy to cause accumulation error in rotation or translation. Compromising between them, the affine model with a moderate DoF, has made a good balance between the aligning accuracy and the global consistency, which gives the most visually satisfactory mosaicking result. Additionally, because of the low flight altitude, the comparatively large-depth-difference ground greatly decreases the aligning precision of the UAV image sequence. In such case, the affine model still shows the similar ability as it does in the aerial image sequence. Conclusively, the affine model has the best comprehensive property to provide a robust initial alignment.

### 3.3. Comprehensive Evaluation on Mosaicking Results

The final mosaicking results of our approach were evaluated in both qualitative and quantitative forms. Firstly, we compare the mosaic images generated by our approach with those composited by a commercial software named PTGui<sup>1</sup> on visual effects. Aiming at comparing the alignment results only, the following seamline detection and tonal

---

<sup>1</sup><http://www.ptgui.com/>

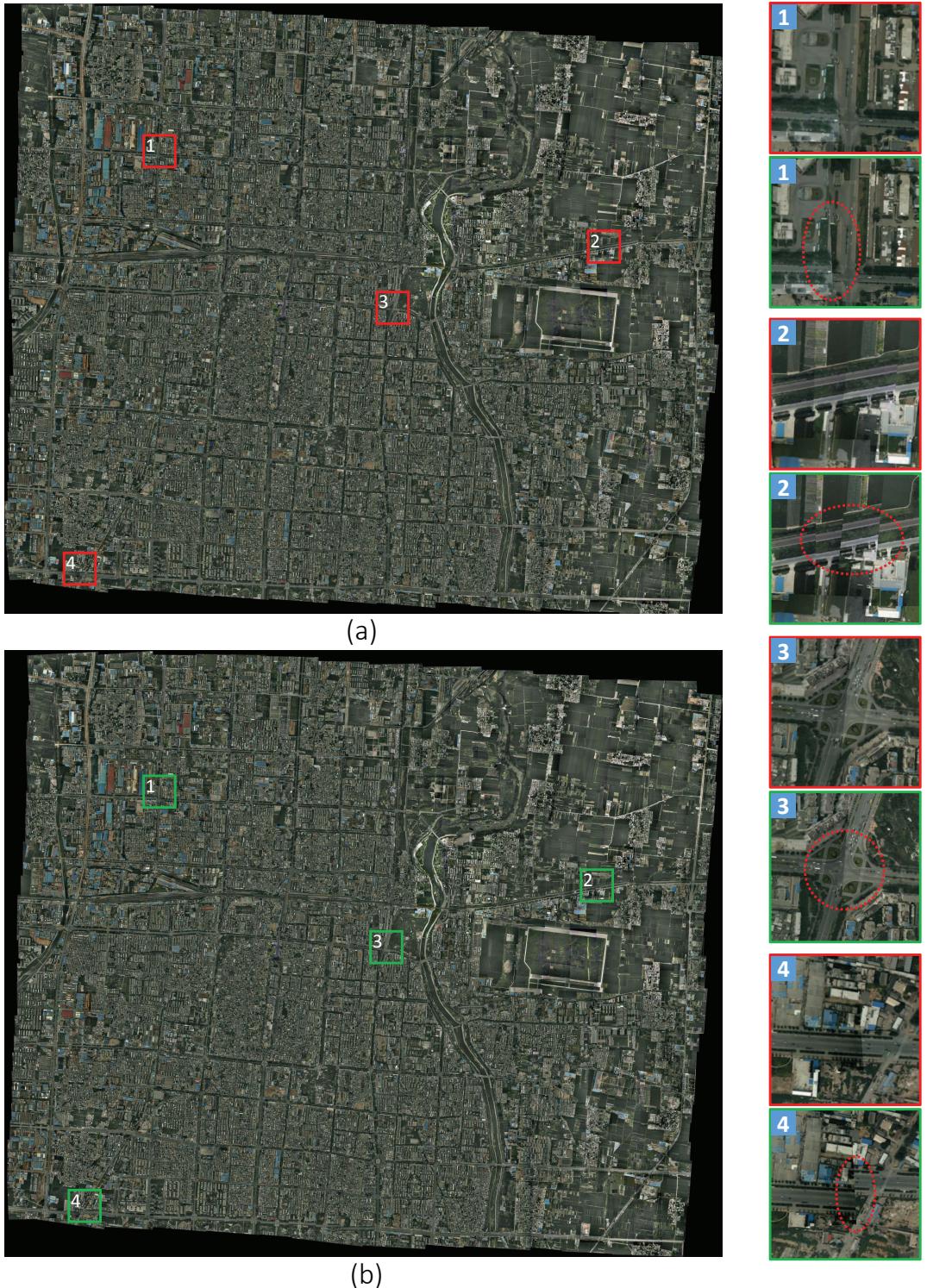


Figure 9: The mosaics composed from the first dataset (744 images) by: (a) our approach and (b) PTGui, respectively. Several typical regions grabbed from the mosaics are enlarged in pairs in the right column.



(a)



(b)



(c)

Figure 10: The mosaics composited from the second dataset (130 images) by: (a) our approach and (c) PTGui, respectively. Several typical regions grabbed from the mosaics are enlarged in pairs in (b).

402 correction were skipped in PTGui and our image stacking order was also made consistent  
 403 with that of PTGui. The comparative results of the first and the second dataset are  
 404 illustrated in Figure 9 and Figure 10, respectively.

405 From the mosaics shown in Figure 9, the two mosaics have similar visual effects as a  
406 whole, both of which take on a pretty good global consistency. However, when it comes to  
407 the local aligning accuracy, our approach has an obvious superiority over PTGui, which  
408 can be observed from some enlarged regions listed in the right column of Figure 9. As for  
409 the UAV data, the large-depth-difference ground makes the assumption of planarity of  
410 the scene weaker, which increases the difficulty to keep the global consistency. A slightly  
411 down-scale tendency in the left part can be found in the mosaicking result of our approach  
412 in Figure 10(a). Since some strong constraints were employed for keeping the scale of  
413 each image consistent, the mosaicking result of PTGui nearly suffered no perspective  
414 distortions, but meanwhile, its alignment precision is destroyed greatly. For a detail  
415 comparison, a serial of enlarged typical regions are listed in the middle line of Figure 10,  
416 which illustrate the better performance of our approach in the aspect of aligning accuracy.

417 Without precision analysis in PTGui, the quantitative evaluation of our approach  
418 was performed in two aspects. As an alignment precision, the registration error of the  
419 initial alignment and the finally global alignment in our approach, were compared in  
420 the form of cumulative probability distribution, as displayed in Figure 8. From the  
421 comparisons, it is easy to find that the aligning precision increases a lot with the help of the  
422 homography refinement, while the global consistency is not affected during the transition  
423 from the affine model to the homography model, as can be observed in Figure 9(a) and  
424 Figure 10(a). This is what we aim at, namely to keep an optimal balance between the  
425 alignment accuracy and the global consistency.

426 In addition, the available poses of the first dataset, recovered by the rigid block  
427 adjustment of photogrammetry field, were used to calculate the homography models  
428 according to the formula in [4] under the assumption of the ground being a plane.  
429 Considering the pesudo-planarity of the ground scene, they are not accurate enough to  
430 be used as the ground truth, but they are qualified to be a reference of global consistency,  
431 since the pose parameters can be regarded as no accumulation error. The recovered  
432 image centroids of the first dataset obtained by our approach and the reference models,  
433 are illustrated in Figure 11. It shows that the two groups of centroids have a similar  
434 distribution form, except for some displacements which average at 5.16 pixels. In fact,

435 as an image mosaicking algorithm based on feature registration, the recovered geometric  
 436 positions are accurate enough to keep the global consistency of a mosaic, which emphasizes  
 437 more on the visual effects than the geometric accuracy. Besides, because of no image  
 438 registration based optimization performed, the pose-based approach has a terrible image  
 439 aligning accuracy with the RMS error of 103.9 pixels, which is much inferior to 1.36 pixels  
 440 obtained by our approach. Therefore, our approach holds a good property of alignment  
 441 accuracy and global consistency in image mosaicking.

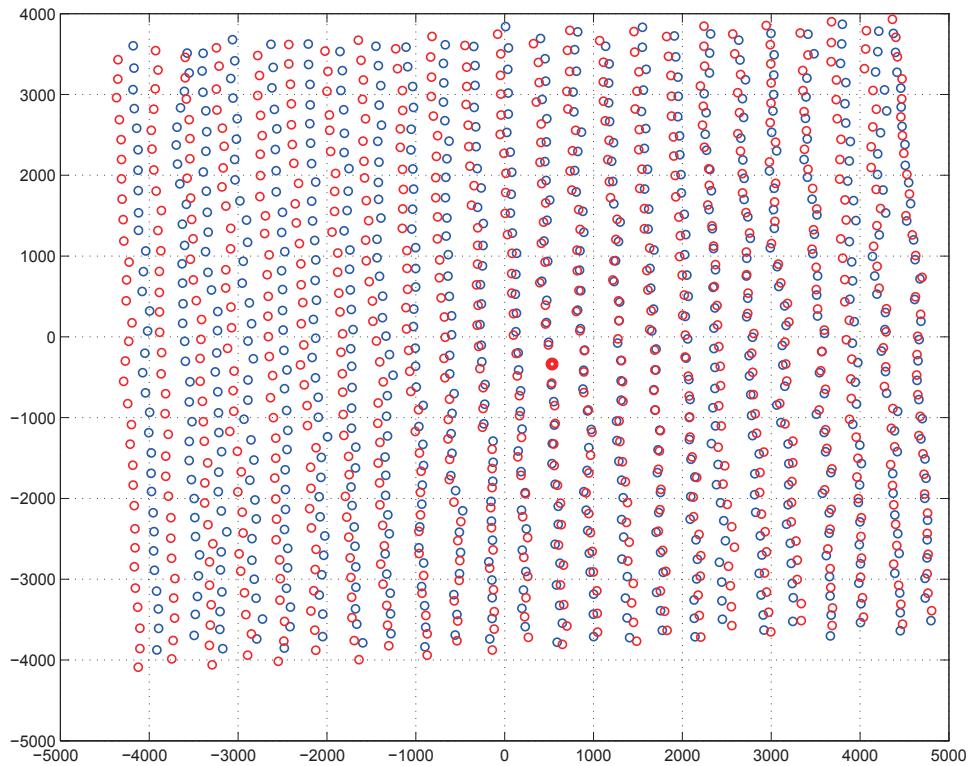


Figure 11: Distributions of image centroids on the mosaic computed by two different approaches. The red circles are the centroids recovered by the proposed approach, and the blue ones represent the result of the pose-based approach. The solid red circle stands for the centroid of the reference image, via which the two groups are strictly superimposed as a base point.

#### 442 4. Conclusion and Future Works

443 In this paper, a topology analysis based generic framework is presented for mosaicking  
 444 sequential images of an approximately planar scene, which contains three steps : topology  
 445 estimation, reference image selection, and global alignment. Specifically, it is adapted

446 to both ordered and unordered image sequences. To estimate the topology robustly, we  
 447 perform the image location and the potential overlapping pairs detection in a collaborative  
 448 way, which makes our approach significantly outperform the state-of-the-art method in  
 449 efficiency. Based on the topological graph, the optimal reference image is found by graph  
 450 analysis and all the images are organized as a spanning tree which gives the reference  
 451 relationships for each image. With the result of topological analysis, we propose a global  
 452 alignment strategy of allowing the continuous transition between the affine model and the  
 453 homography one according to the energy definition, which can keep the optimal balance  
 454 between the global consistency and the aligning accuracy adaptively. The proposed  
 455 framework is tested with several datasets and the experimental results illustrate the  
 456 superiority of our approach. However, the strategy of selecting the reference image, not  
 457 taking the visual angle of image into account, may make a mosaic image of squint angle.  
 458 This problem will be studied in our future work.

## 459 **Acknowledgment**

460 This work was partially supported by the National Natural Science Foundation of  
 461 China (Project No. 41571436), the Hubei Province Science and Technology Support  
 462 Program, China (Project No. 2015BAA027), the National Natural Science Foundation of  
 463 China (Project No. 41271431), and the Jiangsu Province Science and Technology Support  
 464 Program, China (Project No. BE2014866).

## 465 **Appendix A. Optimization Derivation for Model Refinement Under Anti- 466 Perspective Constraint**

467 All the terms in the energy definition in Eq. (8) for model refinement under anti-  
 468 perspective constraint are quadratic, which need to be linearized by the Taylor expansion  
 469 for the iterative optimization. Generally, the first-order Taylor series expansion is accurate  
 470 enough for the optimization problem of quadratic functions.

471 Here, we define the parameter vector of the homography matrix  $\mathbf{H}_i$  as  $\theta_i =$   
 472  $[h_1^i, h_2^i, h_3^i, h_4^i, h_5^i, h_6^i, h_7^i, h_8^i]^\top, i \in [1, n]$ , and the initial value of  $\theta_i$  is defined as  $\bar{\theta}_i =$   
 473  $[\bar{h}_1^i, \bar{h}_2^i, \bar{h}_3^i, \bar{h}_4^i, \bar{h}_5^i, \bar{h}_6^i, \bar{h}_7^i, \bar{h}_8^i]^\top$ . Taking a pair of matching points  $\{\varpi(\mathbf{x}_{ij}^k) = (x, y), \varpi(\mathbf{x}_{ji}^k) =$

<sup>474</sup>  $(x', y')\}$  from  $\mathbf{I}_i$  and  $\mathbf{I}_j$  for example, Eq. (8) can be written as :

$$f_k = \left( \frac{h_1^i x + h_2^i y + h_3^i}{h_7^i x + h_8^i y + 1} - \frac{h_1^j x' + h_2^j y' + h_3^j}{h_7^j x' + h_8^j y' + 1} \right)^2 + \left( \frac{h_4^i x + h_5^i y + h_6^i}{h_7^i x + h_8^i y + 1} - \frac{h_4^j x' + h_5^j y' + h_6^j}{h_7^j x' + h_8^j y' + 1} \right)^2 \\ + \lambda \left[ \left( \frac{h_1^i x + h_2^i y + h_3^i}{h_7^i x + h_8^i y + 1} - x_0 \right)^2 + \left( \frac{h_1^j x' + h_2^j y' + h_3^j}{h_7^j x' + h_8^j y' + 1} - x'_0 \right)^2 \right. \\ \left. + \left( \frac{h_4^i x + h_5^i y + h_6^i}{h_7^i x + h_8^i y + 1} - y_0 \right)^2 + \left( \frac{h_4^j x' + h_5^j y' + h_6^j}{h_7^j x' + h_8^j y' + 1} - y'_0 \right)^2 \right], \quad (\text{A.1})$$

<sup>475</sup> where  $[x_0, y_0]^\top = \varpi(\mathbf{A}_i \mathbf{x}_{ij}^k)$  and  $[x'_0, y'_0]^\top = \varpi(\mathbf{A}_j \mathbf{x}_{ji}^k)$ , are the constant terms which can  
<sup>476</sup> be calculated in advance. Eq. (A.1) is expanded in the form of the first-order Taylor  
<sup>477</sup> series as:

$$f_k \approx \bar{f}_k + \frac{\partial f_k}{\partial h_1^i} dh_1^i + \frac{\partial f_k}{\partial h_2^i} dh_2^i + \frac{\partial f_k}{\partial h_3^i} dh_3^i + \frac{\partial f_k}{\partial h_4^i} dh_4^i + \frac{\partial f_k}{\partial h_5^i} dh_5^i + \frac{\partial f_k}{\partial h_6^i} dh_6^i + \frac{\partial f_k}{\partial h_7^i} dh_7^i + \frac{\partial f_k}{\partial h_8^i} dh_8^i \\ + \frac{\partial f_k}{\partial h_1^j} dh_1^j + \frac{\partial f_k}{\partial h_2^j} dh_2^j + \frac{\partial f_k}{\partial h_3^j} dh_3^j + \frac{\partial f_k}{\partial h_4^j} dh_4^j + \frac{\partial f_k}{\partial h_5^j} dh_5^j + \frac{\partial f_k}{\partial h_6^j} dh_6^j + \frac{\partial f_k}{\partial h_7^j} dh_7^j + \frac{\partial f_k}{\partial h_8^j} dh_8^j, \quad (\text{A.2})$$

<sup>478</sup> where  $\bar{f}_k$  is the values of  $f_k$  when substituting  $\bar{\theta}_i$  and  $\bar{\theta}_j$  into Eq. (A.1).  $d\theta_i =$   
<sup>479</sup>  $[dh_1^i, dh_2^i, dh_3^i, dh_4^i, dh_5^i, dh_6^i, dh_7^i, dh_8^i]^\top$  represents the delta value of  $\theta_i, i \in [1, n]$ . The  
<sup>480</sup> partial derivatives of functions  $f_k$  with respect to  $\theta_i$  and  $\theta_j$  are listed as below:

$$\left\{ \begin{array}{l} \frac{\partial f_k}{\partial h_1^i} = \frac{K_1 x}{h_7^i x + h_8^i y + 1}, \quad \frac{\partial f_k}{\partial h_2^i} = \frac{K_1 y}{h_7^i x + h_8^i y + 1}, \quad \frac{\partial f_k}{\partial h_3^i} = \frac{K_1}{h_7^i x + h_8^i y + 1}, \\ \frac{\partial f_k}{\partial h_4^i} = \frac{K_2 x}{h_7^i x + h_8^i y + 1}, \quad \frac{\partial f_k}{\partial h_5^i} = \frac{K_2 y}{h_7^i x + h_8^i y + 1}, \quad \frac{\partial f_k}{\partial h_6^i} = \frac{K_2}{h_7^i x + h_8^i y + 1}, \\ \frac{\partial f_k}{\partial h_7^i} = \frac{-K_1(h_1^i x + h_2^i y + h_3^i)x}{(h_7^i x + h_8^i y + 1)^2} + \frac{-K_2(h_3^i x + h_4^i y + h_5^i)x}{(h_7^i x + h_8^i y + 1)^2}, \\ \frac{\partial f_k}{\partial h_8^i} = \frac{-K_1(h_1^i x + h_2^i y + h_3^i)y}{(h_7^i x + h_8^i y + 1)^2} + \frac{-K_2(h_3^i x + h_4^i y + h_5^i)y}{(h_7^i x + h_8^i y + 1)^2}, \\ \frac{\partial f_k}{\partial h_1^j} = \frac{K_3 x}{h_7^j x + h_8^j y + 1}, \quad \frac{\partial f_k}{\partial h_2^j} = \frac{K_3 y}{h_7^j x + h_8^j y + 1}, \quad \frac{\partial f_k}{\partial h_3^j} = \frac{K_3}{h_7^j x + h_8^j y + 1}, \\ \frac{\partial f_k}{\partial h_4^j} = \frac{K_4 x}{h_7^j x + h_8^j y + 1}, \quad \frac{\partial f_k}{\partial h_5^j} = \frac{K_4 y}{h_7^j x + h_8^j y + 1}, \quad \frac{\partial f_k}{\partial h_6^j} = \frac{K_4}{h_7^j x + h_8^j y + 1}, \\ \frac{\partial f_k}{\partial h_7^j} = \frac{-K_3(h_1^j x + h_2^j y + h_3^j)x}{(h_7^j x + h_8^j y + 1)^2} + \frac{-K_4(h_3^j x + h_4^j y + h_5^j)x}{(h_7^j x + h_8^j y + 1)^2}, \\ \frac{\partial f_k}{\partial h_8^j} = \frac{-K_3(h_1^j x + h_2^j y + h_3^j)y}{(h_7^j x + h_8^j y + 1)^2} + \frac{-K_4(h_3^j x + h_4^j y + h_5^j)y}{(h_7^j x + h_8^j y + 1)^2}, \end{array} \right.$$

481 where  $K_1$ ,  $K_2$ ,  $K_3$ , and  $K_4$  are computed as:

$$\left\{ \begin{array}{l} K_1 = \frac{2(\bar{h}_1^i x + \bar{h}_2^i y + \bar{h}_3^i)}{\bar{h}_7^i x + \bar{h}_8^i y + 1} - \frac{2(\bar{h}_1^j x' + \bar{h}_2^j y' + \bar{h}_3^j)}{\bar{h}_7^j x' + \bar{h}_8^j y' + 1} + 2\lambda \left( \frac{\bar{h}_1^i x + \bar{h}_2^i y + \bar{h}_3^i}{\bar{h}_7^i x + \bar{h}_8^i y + 1} - x_0 \right), \\ K_2 = \frac{2(\bar{h}_4^i x + \bar{h}_5^i y + \bar{h}_6^i)}{\bar{h}_7^i x + \bar{h}_8^i y + 1} - \frac{2(\bar{h}_4^j x' + \bar{h}_5^j y' + \bar{h}_6^j)}{\bar{h}_7^j x' + \bar{h}_8^j y' + 1} + 2\lambda \left( \frac{\bar{h}_4^i x + \bar{h}_5^i y + \bar{h}_6^i}{\bar{h}_7^i x + \bar{h}_8^i y + 1} - y_0 \right), \\ K_3 = -\frac{2(\bar{h}_1^i x + \bar{h}_2^i y + \bar{h}_3^i)}{\bar{h}_7^i x + \bar{h}_8^i y + 1} + \frac{2(\bar{h}_1^j x' + \bar{h}_2^j y' + \bar{h}_3^j)}{\bar{h}_7^j x' + \bar{h}_8^j y' + 1} + 2\lambda \left( \frac{\bar{h}_1^j x' + \bar{h}_2^j y' + \bar{h}_3^j}{\bar{h}_7^j x' + \bar{h}_8^j y' + 1} - x'_0 \right), \\ K_4 = -\frac{2(\bar{h}_4^i x + \bar{h}_5^i y + \bar{h}_6^i)}{\bar{h}_7^i x + \bar{h}_8^i y + 1} + \frac{2(\bar{h}_4^j x' + \bar{h}_5^j y' + \bar{h}_6^j)}{\bar{h}_7^j x' + \bar{h}_8^j y' + 1} + 2\lambda \left( \frac{\bar{h}_4^j x' + \bar{h}_5^j y' + \bar{h}_6^j}{\bar{h}_7^j x' + \bar{h}_8^j y' + 1} - y'_0 \right). \end{array} \right.$$

482 For the convenience of descriptions in the following, the matrix form of Eq. (A.2) are  
483 written as the standard equation of the Least Square optimization:

$$[v_k] = \left[ \begin{array}{cccccccccc} \dots & \frac{\partial f_k}{\partial h_1^i} & \frac{\partial f_k}{\partial h_2^i} & \frac{\partial f_k}{\partial h_3^i} & \frac{\partial f_k}{\partial h_4^i} & \frac{\partial f_k}{\partial h_5^i} & \frac{\partial f_k}{\partial h_6^i} & \frac{\partial f_k}{\partial h_7^i} & \frac{\partial f_k}{\partial h_8^i} & \dots \\ \dots & \frac{\partial f_i}{\partial h_1^j} & \frac{\partial f_i}{\partial h_2^j} & \frac{\partial f_i}{\partial h_3^j} & \frac{\partial f_i}{\partial h_4^j} & \frac{\partial f_i}{\partial h_5^j} & \frac{\partial f_i}{\partial h_6^j} & \frac{\partial f_i}{\partial h_7^j} & \frac{\partial f_i}{\partial h_8^j} & \dots \end{array} \right] \begin{bmatrix} d\theta_i \\ \vdots \\ d\theta_j \\ \vdots \end{bmatrix} - [-\bar{f}_k].$$

484 The above equation is expressed with the corresponding matrix labels as:

$$\mathbf{V}^k = \mathbf{J}^k \mathbf{X} - \mathbf{L}^k, \quad (\text{A.3})$$

485 where the dots in the Jacobi matrix  $\mathbf{J}^k$  represent a series of zeros, and the dots in  $\mathbf{X}$   
486 indicate the other unknown parameters in  $\{d\theta_i\}_{i=1}^n$ .  $\mathbf{V}^k$  is the residual error of a pair of  
487 matching points. Hereafter, we name  $\mathbf{J}^k$  and  $\mathbf{L}^k$  as the coefficient matrix and the constant  
488 matrix, respectively.

489 As can be seen, a pair of matching points from two images provides an equation with  
490 16 unknown parameters. Supposing that  $n$  images have  $m$  pairs of overlapping relations  
491 and there are  $s$  matching points of each image pair in average, then we obtain a Jacobi  
492 matrix with the size of  $m \times s$  rows and  $8 \times m$  columns and a constant matrix with the size  
493 of  $m \times s$  rows and 1 column. In each iteration,  $\mathbf{J}_{ms \times 8n}$  and  $\mathbf{L}_{ms \times 1}$  have to be recalculated  
494 and the corresponding solution vector  $\mathbf{X}_{8n \times 1} = [d\theta_1^\top, \dots, d\theta_n^\top]^\top$  can be solved with the  
495 following equation:

$$\mathbf{X}_{8n \times 1} = (\mathbf{J}_{ms \times 8n}^\top \mathbf{J}_{ms \times 8n})^{-1} (\mathbf{J}_{ms \times 8n}^\top \mathbf{L}_{ms \times 1}). \quad (\text{A.4})$$

496 The initial solution of  $\{\theta_i\}_{i=1}^n$  for next iteration is updated by adding up  $\mathbf{X}_{8n \times 1}$  and the  
 497 initial solution used in this iteration. As the iteration goes, the updated solution will  
 498 converge to the optimal solution gradually unless the initial solution provided at the very  
 499 beginning is not accurate enough. However, when the amount of images is large, the size  
 500 of the Jacobi matrix will be very huge and makes a challenge to the memory of computer.

501 In fact, we can calculate  $\{\theta_i\}_{i=1}^n$  directly if  $\mathbf{TJ}_{8n \times 8n} = \mathbf{J}_{ms \times 8n}^\top \mathbf{J}_{ms \times 8n}$  and  $\mathbf{TL}_{8n \times 1} =$   
 502  $\mathbf{J}_{ms \times 8n}^\top \mathbf{L}_{ms \times 1}$  have been obtained. So, to reduce the required memory space and the  
 503 computation time, we manage to compute  $\mathbf{TJ}_{8n \times 8n}$  and  $\mathbf{TL}_{8n \times 1}$  by adding up the matrix  
 504  $\mathbf{J}^{i^\top} \mathbf{J}^i$  and the matrix  $\mathbf{J}^{i^\top} \mathbf{L}^i$  calculated from each pair of matching points, instead of  
 505 building the large  $\mathbf{J}$  and  $\mathbf{L}$  beforehand. The improved computation formula is defined as:

$$\begin{cases} \mathbf{TJ}_{8n \times 8n} = \sum_{i=1}^{ms} \mathbf{J}_{1 \times 8n}^{i^\top} \mathbf{J}_{1 \times 8n}^i, \\ \mathbf{TL}_{8n \times 1} = \sum_{i=1}^{ms} \mathbf{J}_{1 \times 8n}^{i^\top} \mathbf{L}_{2 \times 1}^i. \end{cases} \quad (\text{A.5})$$

506 Then, the solution can be obtained in this way as:

$$\mathbf{X}_{8n \times 1} = \mathbf{TJ}_{8n \times 8n}^{-1} \mathbf{TL}_{8n \times 1}. \quad (\text{A.6})$$

507 Additionally, considering the sparsity of  $\mathbf{J}^i$ , the computation of the matrix multiplication  
 508 in Eq. (A.5) can be improved further in the complexity of both time and space.

## 509 References

- 510 [1] E. Zagrouba, W. Barhoumi, S. Amri, An efficient image-mosaicing method based on  
 511 multifeature matching, *Machine Vision and Applications* 20 (3) (2009) 139–162.
- 512 [2] J. Chen, H. Feng, K. Pan, Z. Xu, Q. Li, An optimization method for registration and  
 513 mosaicking of remote sensing images, *Optik - International Journal for Light and Electron  
 514 Optics* 125 (2) (2014) 697–703.
- 515 [3] B. Zitova, J. Flusser, Image registration methods: a survey, *Image and Vision Computing*  
 516 21 (11) (2003) 977–1000.
- 517 [4] L. Kang, L. Wu, Y. Wei, B. Yang, H. Song, A highly accurate dense approach for  
 518 homography estimation using modified differential evolution, *Engineering Applications of  
 519 Artificial Intelligence* 31 (4) (2014) 68–77.

- 520 [5] Z. Wang, Y. Chen, Z. Zhu, W. Zhao, An automatic panoramic image mosaic method based  
521 on graph model, *Multimedia Tools and Applications* 75 (5) (2015) 2725–2740.
- 522 [6] N. R. Gracias, S. Van Der Zwaan, A. Bernardino, S.-V. Jos, Mosaic-based navigation  
523 for autonomous underwater vehicles, *IEEE Journal of Oceanic Engineering* 28 (4) (2003)  
524 609–624.
- 525 [7] Y. Xu, J. Ou, H. He, X. Zhang, J. Mills, Mosaicking of unmanned aerial vehicle imagery  
526 in the absence of camera poses, *2016 8* (3) (2016) 204.
- 527 [8] Y. He, R. Chung, Image mosaicking for polyhedral scene and in particular singly visible  
528 surfaces, *Pattern Recognition* 41 (3) (2008) 1200–1213.
- 529 [9] J. Zaragoza, T.-J. Chin, M. Brown, D. Suter, As-projective-as-possible image stitching  
530 with moving DLT, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36 (7)  
531 (2014) 1285–1298.
- 532 [10] C.-H. Chang, Y. Sato, Y.-Y. Chuang, Shape-preserving half-projective warps for image  
533 stitching, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp.  
534 3254–3261.
- 535 [11] D. Patidar, A. Jain, Automatic image mosaicing: an approach based on FFT, *International  
536 Journal of Scientific Engineering and Technology* 1 (1) (2011) 01–04.
- 537 [12] S. Ghannam, A. L. Abbott, Cross correlation versus mutual information for image  
538 mosaicing, *International Journal of Advanced Computer Science and Applications* 4 (11)  
539 (2013) 94–102.
- 540 [13] A. Elibol, R. Garcia, O. Delaunoy, N. Gracias, *Efficient Topology Estimation for Large  
541 Scale Optical Mapping*, Springer, 2013, Ch. A New Global Alignment Method for Feature  
542 Based Image Mosaicing, pp. 25–39.
- 543 [14] S. Ali, C. Daul, E. Galbrun, Guillemin, Anisotropic motion estimation on edge preserving  
544 Riesz wavelets for robust video mosaicing, *Pattern Recognition* 51 (2016) 425–442.
- 545 [15] Y. Chen, J. Sun, G. Wang, Minimizing geometric distance by iterative linear optimization,  
546 in: *IEEE International Conference on Pattern Recognition*, Vol. 29, 2010, pp. 1–4.

- 547 [16] W. Mou, H. Wang, G. Seet, L. Zhou, Robust homography estimation based on nonlinear  
548 least squares optimization, Mathematical Problems in Engineering 2014 (1) (2013) 372–377.
- 549 [17] L. Zhang, Y. Li, J. Zhang, Y. Hu, Homography estimation in omnidirectional vision under  
550 the  $L_\infty$ -norm, in: IEEE International Conference on Robotics and Biomimetics, 2010, pp.  
551 1468–1473.
- 552 [18] B. Triggs, P. F. McLauchlan, R. I. Hartley, A. W. Fitzgibbon, Vision Algorithms: Theory  
553 and Practice, Lecture Notes in Computer Science, Springer, 2000, Ch. Bundle adjustment and  
554 modern synthesis, pp. 298–372.
- 555 [19] K. Konolige, Sparse sparse bundle adjustment, in: British Machine Vision Conference,  
556 2010, pp. 1–10.
- 557 [20] M. Li, D. Li, D. Fan, A study on automatic UAV image mosaic method for paroxysmal  
558 disaster, in: Proceedings of the International Society of Photogrammetry and Remote  
559 Sensing Congress, 2012.
- 560 [21] C. Xing, J. Wang, Y. Xu, A robust method for mosaicking sequence images obtained from  
561 UAV, in: International Conference on Information Engineering and Computer Science  
562 (ICIECS), 2010.
- 563 [22] F. Caballero, L. Merino, J. Ferruz, A. Ollero, Homography based Kalman filter for mosaic  
564 building. applications to UAV position estimation, in: IEEE International Conference on  
565 Robotics and Automation, 2007, pp. 2004–2009.
- 566 [23] A. Y. Taygun Kekec, M. Unel, A new approach to real-time mosaicing of aerial images,  
567 Robotics and Autonomous Systems 62 (12) (2014) 1755–1767.
- 568 [24] A. Elibol, R. Garcia R Elibol A, Gracias Na, O. Delaunoy, N. Gracias, A new global  
569 alignment method for feature based image mosaicing, Advances in Visual Computing,  
570 Lecture Notes in Computer Science 5359 (7) (2008) 257–266.
- 571 [25] M. Xia, J. Yao, L. Li, X. Lu, Globally consistent alignment for mosaicking aerial images,  
572 in: IEEE International Conference on Image Processing, 2015.
- 573 [26] E.-Y. Kang, I. Cohen, G. Medioni, A graph-based global registration for 2D mosaics, in:  
574 International Conference on Pattern Recognition, Vol. 1, 2000, pp. 257–260.

- 575 [27] R. Marzotto, A. Fusiello, V. Murino, High resolution video mosaicing with global alignment,  
576 in: IEEE Conference on Computer Vision and Pattern Recognition, Vol. 1, 2004, pp. 692–  
577 698.
- 578 [28] A. Elibol, N. Gracias, R. Garcia, Fast topology estimation for image mosaicing using  
579 adaptive information thresholding, *Robotics and Autonomous systems* 61 (2) (2013) 125–  
580 136.
- 581 [29] R. Szeliski, Image alignment and stitching: A tutorial, *Foundations and Trends in*  
582 *Computer Graphics and Vision* 2 (1) (2006) 1–104.
- 583 [30] T. E. Choe, I. Cohen, M. Lee, G. Medioni, Optimal global mosaic generation from retinal  
584 images, in: IEEE International Conference on Pattern Recognition, Vol. 3, 2006, pp. 681–  
585 684.
- 586 [31] B. Bollobas, *Modern graph theory*, Springer Science & Business Media, 2013.
- 587 [32] R. L. Graham, P. Hell, On the history of the minimum spanning tree problem, *Annals of*  
588 *the History of Computing* 7 (1) (1985) 43–57.
- 589 [33] T. T. Bay, Herbert, L. V. Gool, Surf: Speeded up robust features, in: IEEE European  
590 Conference on Computer Vision, 2006.
- 591 [34] R. W. Floyd, Algorithm 97: shortest path, *Communications of the ACM* 5 (6) (1962)  
592 345–345.
- 593 [35] T. H. Cormen, *Introduction to algorithms*, MIT press, 2009.
- 594 [36] R. I. Hartley, In defense of the eight-point algorithm, *IEEE Transactions on Pattern*  
595 *Analysis and Machine Intelligence* 19 (6) (1997) 580–593.
- 596 [37] P. Torr, A. Zisserman, MLESAC: A new robust estimator with application to estimating  
597 image geometry, *Computer Vision and Image Understanding* 78 (1) (2000) 138–156.
- 598 [38] M. I. A. Lourakis, Sparse non-linear least squares optimization for geometric vision, in:  
599 Computer Vision–ECCV 2010, Vol. 6312, 2010, pp. 43–56.