

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

Over the past decade, many image mosaicking methods have been proposed in robotic mapping and remote sensing applications. However, most of these methods mainly focus on the optimizing problem of minimizing the image alignment error, which can't necessarily guarantee the global consistency of the mosaicking result, especially for the case of wide-range pseudo-planar scenes prone to suffering from severe perspective distortion. 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 perspective distortion meanwhile preserving 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 unordered image sequence), along which potential overlapping pairs are incrementally detected according to the gradually recovered geometrical positions and orientations. Secondly, all the sequential images are organized as a spinning tree through applying a graphic algorithm on the topological graph, so as to find the optimal reference image which minimizes the total number of error propagation. Thirdly, we perform the global consistent alignment with the topology analysis in an ingenious strategy that images are initially aligned by groups via the robust affine model, followed by the model

*Corresponding author.

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

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

refinement under the anti-perspective constraint, through which the optimal balance between aligning precision and global consistency can be achieved. Finally, experimental results on several challenging aerial image sets sufficiently illustrate the validity of the proposed approach.

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

1. Introduction

Owing to the rapid developments in obtaining optical image data from areas beyond human reach, there is a high demand from different research and engineering fields for creating large range mosaicked images. In fact, image mosaicking is a procedure that merges two or more images with overlapping areas into a single composite image as seamless as possible in both geometry and color tone. The critical first 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]. Because of the limitation of motionless position, two or multiple images captured from the same camera center and toward different orientations are mainly used to make a ground panorama with an omniscient point of view [5]. On the contrary, mosaicking images of a 3D planar scene permits the camera moving freely, which is popular in robotic mapping and remote sensing applications [6, 7]. Recently, some mosaicking methods not limited to this 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 known as planar image mosaicking. Under the challenge of both pseudo-plane and accumulation error, a lot of related studies have been presented in the literature of the last decade. However, the performance considering 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

24 cost, the area-based approaches are seldom used in the mosaicking missions of large
25 scale [14]. As for the planar mosaicking problem, such as aerial image mosaicking, the
26 feature-based approaches are usually applied to recover the homography model between
27 images [15, 16, 17] due to the fact that the ground scene can be regarded as an approximate
28 plane observed from the aerial photographic camera. To improve the mosaicking result,
29 many optimization algorithms have been proposed to achieve a global alignment. A
30 typical global optimization method is "Bundle Adjustment" [18, 19], which aims at finding
31 an optimal solution minimizing the total reprojection error [20]. To provide a good initial
32 solution for global optimization, Xing et al. [21] proposed to first apply the Extended
33 Kalman Filter [22] onto the local area, and then refine all the parameters globally. To
34 avoid the non-linear optimization, Kekec et al. [23] employed the affine model to optimize
35 the initial alignment made by the homography model in the global optimization. Some
36 methods [24, 25] utilized the topological structure information of images to achieve a
37 global registration. To prevent image suffering down-scaling effect, Elibol et al. [26]
38 proposed to optimize point positions in the mosaicked frame and the alignment model in
39 an alternate iteration scheme.

40 All those methods concentrating on the optimizing strategy seeking for an alignment
41 with the least registration error can usually composite a satisfied mosaic image from
42 several or dozens of images. However, sequential images taken from a wide-range area
43 can always make the global consistency inaccessible for them, because in the case of
44 pseudo-plane violating the strict geometric model, the least-registration-error principle
45 is prone to causing a severe accumulation of perspective distortions. To release this
46 problem, Caballero et al. [22] proposed to use the hierarchical models according to the
47 alignment quality of images, where the model with less degree of freedom (DoF) is used
48 for images with bigger parallax. The essence of this method is to make a trade-off
49 between improving aligning precision and resisting perspective distortion. In fact, a more
50 reasonable solution is to allow continuous transition between aligning models according
51 the predefined constraint, instead of regarding the model selection as a binary problem.
52 This idea has been detailedly investigated in our previous work [27].

53 As to the large-scale mosaicking problem, utilizing the topology among images is

54 another effective way to improve the mosaicking result. On the one hand, the potential
55 overlapping relations in topology contribute on the global alignment greatly by providing
56 lots of joint constraints, on the other hand, based on the topological graph, some graphic
57 algorithms can be applied to optimize aligning strategy and reduce error propagation. To
58 estimate the topology efficiently, Elibol et al. [28] used the low-cost tentative matching
59 combined with the Minimum Spanning Tree (MST) solution to detect overlapping
60 relations in an iterative scheme and decide when to update the topological estimation
61 via information-theory principles. The algorithm is efficient as a whole, but the strategy
62 of detecting potential overlapping pairs is not efficient enough, because the detection and
63 the alignment are divided into two independent steps in each iteration, which induce
64 many invalid matching attempts. As for the selection of the reference image, Richard et
65 al. [29] stated that a reasonable choice is the most central image geometrically. This idea
66 is obviously reasonable due to the fact that the central image usually has the shortest
67 distance to all other images on average. However, they didn't give any solution about
68 how to find such an image. To solve this problem, Choe et al. [30] applied a graphic
69 algorithm to select the reference image with the lowest cumulative registration error, but
70 the registration error between each image pair have to be calculated in advance.

71 In this paper, for mosaicking images taken from a wide-range approximately planar
72 scene, we propose to achieve a visually satisfactory mosaic image with both accurate
73 alignment and global consistency through two technical means: (1) utilizing topology
74 analysis to strengthen registration constraints and reduce error propagation; (2) adopting
75 the alignment strategy of allowing continuous transition between different aligning
76 models, to adaptively keep the optimal balance between alignment accuracy and
77 global consistency (i.e., no obviously perspective distortion). Firstly, we initialize an
78 approximate similarity matrix for image pairs in a fast way, which is combined with
79 the Minimum Spanning Tree (MST) to find the main chain for an unordered image
80 sequence. Then, other potential overlaps are detected incrementally with the gradually
81 recovered geometric positions along the main chain. Because of the synchronism of
82 overlap detection and image location, our proposed topology estimation strategy is more
83 efficient than the method used in [28]. Secondly, all the sequential images are organized

84 as a spinning tree through the classical Floyd-Warshall algorithm, so as to find the
85 optimal reference image with the least cascading times when projecting other images
86 to the reference plane. Obviously, such a reference image benefits in reducing error
87 accumulation. Finally, a globally consistent alignment strategy is proposed to align
88 images into a common coordinate system, which combines the affine model with the
89 homography model effectively. The initial alignment is made by the robust affine model
90 by groups and the globally homography refinement is followed under the anti-perspective
91 constraint, to improve the alignment accuracy on the premise of global consistency
92 not affected. Our proposed approach was sufficiently examined through several groups
93 of experiments on two challenging aerial image datasets and the performances were
94 comprehensively evaluated by comparing with the state-of-the-art algorithm and a famous
95 commercial software.

96 The remainder of this paper is organized as follows. The proposed framework is
97 detailed in Section 2, which is comprised of topology estimation, selection of reference
98 image, and global alignment. Experimental results are provided in Section 3 followed by
99 the conclusions are drawn and future works are provided in Section 4.

100 2. Our Approach

101 Aiming at achieving the mosaicking result with both accurate alignment and global
102 consistency, we propose a generic framework for globally consistent alignment of images
103 captured from an approximately planar scene as shown in Figure 1, which is composed of
104 three modules: topology estimation, selection of reference image, and global alignment.
105 First, the sequential images are inputed for topology estimation, through which the
106 obtained topological graph and matching results are used to search the optimal reference
107 image via graph algorithm and to provide feature correspondences for the global
108 alignment, respectively. Finally, according to the reorganized aligning hierarchy, all the
109 images are aligned by a specially designed double-model under the global optimization
110 framework. Due to the versatile topology estimation, the proposed framework is suitable
111 for both time-consecutive image sets and unordered image sets.

112 For the description convenience in the following, the frequently used notations in this

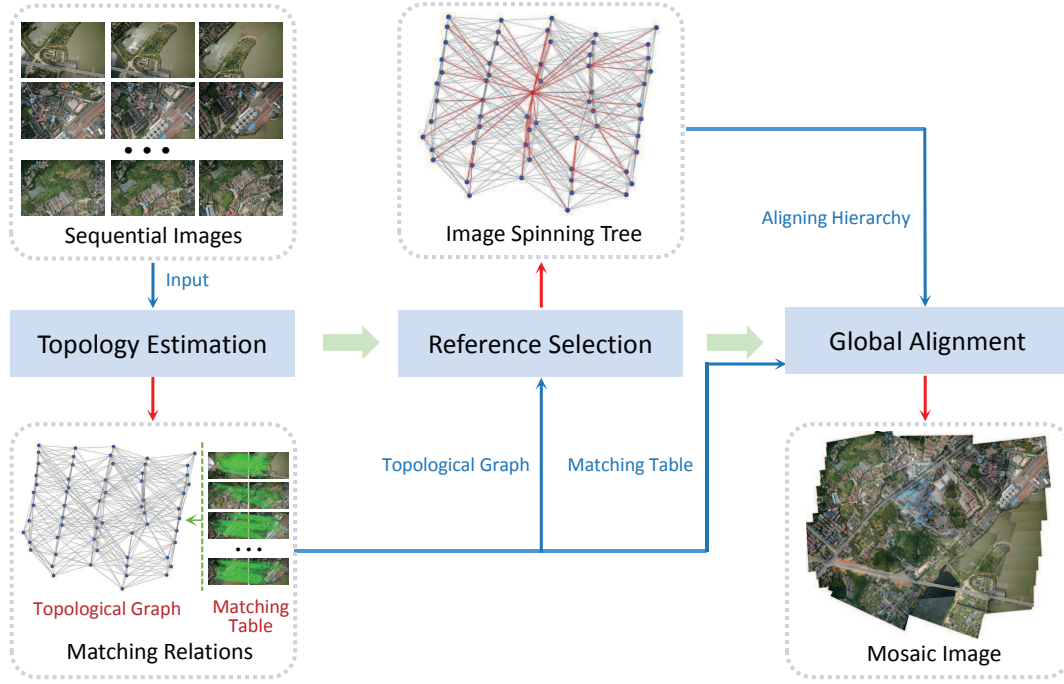


Figure 1: The flowchart of our proposed framework for globally consistent alignment of images. The blue and red thin arrows denote the input and output of each processing module, respectively, and the wide green arrows indicate the execution sequence.

113 paper are summarized below :

- 114 • \mathbf{I}_i - the i -th image in the sequential images.
- 115 • \mathbf{A}_i - the 3×3 affine transformation matrix (6DoF) relating \mathbf{I}_i to the reference frame.
- 116 • \mathbf{H}_i - the 3×3 homography transformation matrix (8DoF) relating \mathbf{I}_i to the reference
- 117 frame.
- 118 • $\mathbf{x} = [x, y, w]^\top$ - the homogeneous coordinate of a feature point.
- 119 • $\mathbf{x}_{i,j}^k$ - the 2D coordinate of the k -th matched feature in \mathbf{I}_i corresponding to the k -th
- 120 matched feature $\mathbf{x}_{j,i}^k$ in \mathbf{I}_j .
- 121 • $M_{i,j}$ - the total number of matches between \mathbf{I}_i and \mathbf{I}_j .
- 122 • $\varpi(\mathbf{x}) = [x/w, y/w]^\top$ - the function transforming the homogeneous coordinate of a
- 123 2D point into the non-homogeneous coordinate.

124 2.1. Fast Topology Estimation

125 The image topology of the surveyed area is usually represented by a graph where an
126 image stands for a node and the overlapping relationship between image pair is denoted by
127 an edge or a link. Topology estimation means to find the existing overlapping relationships
128 among all images. In this section, we try to find all the potential overlapping image pairs
129 by utilizing the gradually recovered geometric positions of images in the time-consecutive
130 order on the mosaicking plane, instead of simply doing matching attempts. As for an
131 unordered image sequence, finding a main chain connecting all images can make the
132 problem the same as that of the time-consecutive image sequence. Therefore, an efficient
133 strategy can be proposed to find the complete topology with the minimum number of
134 image matching attempts.

135 2.1.1. Finding Main Chain with Most Reliability

136 For a sequence of n images, the main chain consisting of $(n - 1)$ edges connects all the
137 nodes/images in the graph. More strictly, it is defined as a spanning tree of an undirected
138 graph in graphic theory [31, 32]. Obviously, there is no need to find a main chain for the
139 time-consecutive image sequence due to that their time-consecutive links have implied a
140 main chain already. That's to say, this step is mainly set for the case of finding the image
141 topology of an unordered image set.

142 Given an unordered image set, we have to measure the similarities between image
143 pairs in advance of finding a main chain. Here, the similarity measurement is intended
144 to be computed in an approximate but efficient way. To achieve this goal, for each
145 image, we select a subset of SURF features extracted from it, and the similarity between
146 image pair is defined as the number of candidate point matches whose descriptor vector
147 distances are less than some given distance. Specially, to increase the corresponding
148 probability, the subsets are generated by selecting features extracted from the same scale
149 layer in the SURF detector, instead of sampling randomly. In our experiments, the
150 features from the second scale layer of the total four octaves were selected as the subset
151 representing each image, which hold a stable ratio of $22 \pm 3\%$ almost for all kinds of
152 images. The computational cost of this similarity measure is comparatively low, since it
153 mainly involves computing the distances between a small set of descriptor vectors. Over

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

159 Although the similarity table built by this way is not reliable, it is qualified to provide
160 the initial similarity information just for finding a main chain under an iterative scheme.
161 Based on the similarity matrix, the reciprocals of those non-zero similarity values are set
162 as the weights for the edges of the graph, i.e., $W(i, j) = \frac{1}{\mathbf{S}(i, j)}$. Given such a graph, we try
163 to select a linkage path that connects all the nodes with the highest total reliability, i.e.,
164 the lowest sum of weights. This idea is effectively implemented in the following two-step
165 iterative scheme.

166 **Maximum Reliability:** This is realized by finding the Minimum Spanning Tree
167 (MST) of the current weighted graph. The MST is a spanning tree whose edges have the
168 minimum total weight in all the spanning trees of the graph. So, the MST represents the
169 connected tree composed of the most similar image pairs.

170 **Check Connectivity:** The algorithm tries to match all the image pairs in the MST.
171 If all the matching attempts succeed completely, the MST is the targeted main chain and
172 the iteration is terminated. On the contrary, when there exists any image pair failed to
173 be matched, we have to modify the weights of the graph where the weights of successfully
174 matched pairs are set as zero while the weights of matching-failed image pairs are set as
175 an infinite value, then it turns to the next iteration.

176 To examine the difference of the main chains from a time-consecutive image sequence
177 and an unordered image sequence, a subset of the first dataset described in Section 3 was
178 selected to demonstrate the results of topology estimation, which had been performed
179 in both the time-consecutive mode and the unordered one, respectively, as shown in
180 Figure 2. Apart from the difference of the main chains, the topologies estimated in these
181 two different modes are almost the same, which are compared quantitatively in Table 1.

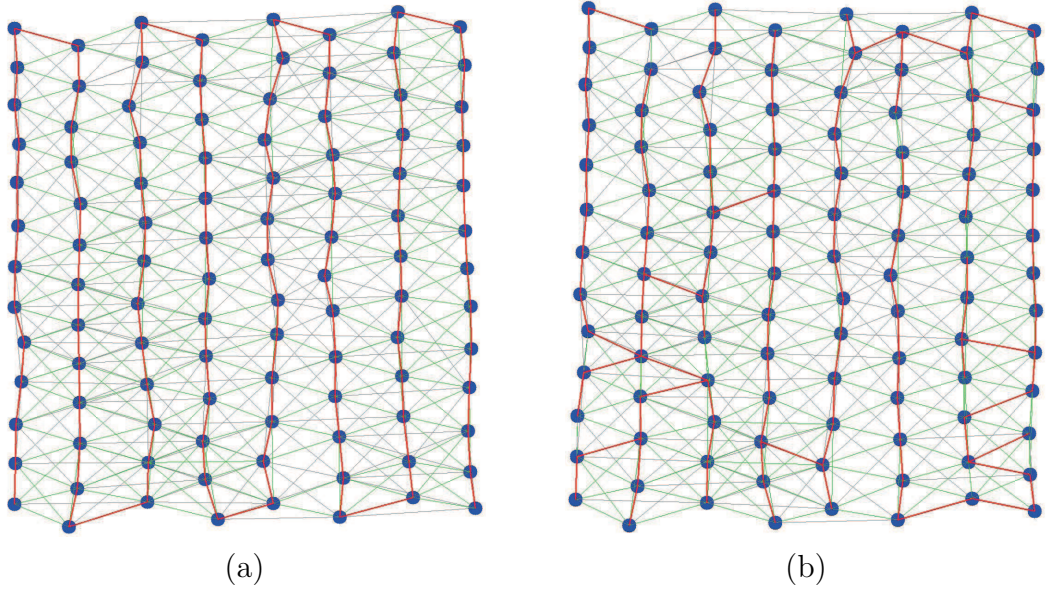


Figure 2: The estimated topologies of an image sequence (104 images) in the time-consecutive mode and the unordered one, respectively: (a) The topology estimated in the time-consecutive mode, where the red edges represent the prior main chain in the time-consecutive order and green edges indicate the numbers of matched features between image pairs are over 100 while gray ones indicate the numbers are less than 100; (b) The topology estimated in the unordered mode, where the red edges represent the main chain linked by the proposed iterative scheme and other edges have the same meanings as in (a).

182 2.1.2. Detecting Potential Overlapping Pairs

183 After having obtained the overlapping relationships along the main chain, we move on
 184 to detect other potential overlapping pairs for a more complete topology. As mentioned
 185 in the beginning of Section 2.1, we can recover the comparative geometric positions of
 186 sequential images in a common coordinate system according to the main chain. Based
 187 on the geometric information, the potential overlapping pairs can be detected easily.
 188 Therefore, there are two problems to be solved: 1) how to recover the comparative
 189 geometric positions with the main chain information; 2) how to detect the potential
 190 overlapping pairs based on the geometric relationships. In the proposed method, these
 191 two problems are solved in a collaborative way, instead of an independent way.

192 Firstly, we temporarily select a reference image as the mosaicking plane through
 193 applying the algorithm detailed in Section 2.2 on the main chain. To recover the
 194 comparative geometric positions, we employ the affine model to align images into the
 195 mosaicking plane, which is robust in locating the centroids of images. Compared with

Algorithm 1 Detecting potential overlapping pairs

Input: The image set $\mathcal{I} = \{\mathbf{I}_i\}_{i=1}^n$ arranged in some order.

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  $\mathcal{P}$ 
```

196 the affine model, the classic homography model is prone to suffering from the perspective
197 distortion, and the 2D rigid model tends to make a bending trajectory because of error
198 accumulation, which is validated in Section 3.2. Specially, to improve the reliability
199 of the image locations, the images on the main chain are aligned starting from the
200 reference image one by one. As the images being located gradually, the potential
201 overlapping relationships around the newly located image are detected and would be
202 used for optimizing the position of this newly aligned image in the following. This
203 strategy makes a significant contribution to improving the accuracy of the recovered
204 geometric positions, because the simultaneously detected overlapping pairs can provide
205 extra favorable constraints for aligning images. Given a newly aligned image \mathbf{I}_i , it checks
206 whether there is any overlap with all the previously aligned image set $\widehat{\mathcal{I}} = \{\mathbf{I}_j\}_{j=1}^m$. For
207 an image pair \mathbf{I}_i and \mathbf{I}_j , the overlap detection is performed by calculating the distance

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

208 between their centroids as follows:

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

209 where c_i , c_j , d_i and d_j are the centroids and the diameters of the minimum boundary
210 circles of the projection onto the mosaicking plane of \mathbf{I}_i and \mathbf{I}_j , respectively. If $\delta_{ij} > 1$,
211 there is no overlap. Otherwise, there may exist an overlap between \mathbf{I}_i and \mathbf{I}_j , and we
212 try to match them for verification. Of course, if the matching between \mathbf{I}_i and \mathbf{I}_j has
213 been attempted during finding the main chain, there is no need to repeat the matching
214 attempt again. The whole sketch procedure of our proposed topology estimation approach
215 is described in Algorithm 1. When all the overlapping pairs are obtained, we redefine the
216 similarity matrix as the final topological representation. The original similarity matrix
217 is reset as a zero matrix firstly, and the value of $\mathbf{S}(i, j)$ is replaced with the number of
218 matched points only if \mathbf{I}_i and \mathbf{I}_j have been matched successfully.

219 It should be noted that the major computation cost of the topology estimation is
220 feature matching between images, as listed in the fourth column of Table 1. Because
221 there is no global optimization or iterative detection, the image alignment and potential
222 overlapping detection have a relatively low computation cost. Besides, as for an unordered
223 image set, the initialization of the similarity matrix occupies the majority of the rest
224 computation actually. That’s to say, the topology estimation before image mosaicking
225 is well worthy, which is fundamental to the following process while adds nearly no extra
226 computation load except for the necessary feature matching.

227 *2.2. Optimal Reference Image Selection*

228 As we known, the images alignment is realized through warping each image into the
 229 mosaicking plane which is always set by selecting one of the sequential images (named
 230 as the reference image). An image without direct overlap to the reference image has
 231 to be projected to the mosaicking plane by cascading a series of relative transformation
 232 models between other intermediate images. Obviously, less intermediate images used for
 233 cascading makes less error accumulation. In fact, there may exist more than one path
 234 with the same cascading numbers from an image to another. Considering each cascading
 235 implies a different error, we would rather to select the path with the least accumulation
 236 error. In terms of this, the optimal reference image should give the lowest sum of
 237 accumulation errors from all the other images to the reference image plane. To address
 238 this problem, we construct an undirected and weighted graph based on the estimated
 239 topology in Section 2.1. According to the similarity matrix obtained by the topology
 240 estimation, those image pairs with non-zero values of similarities are linked with edges.
 241 As far as the weight (or cost) of an edge concerned, there are two kinds of settings in
 242 the existing literature: the reciprocal of the number of matched features [28] and the
 243 registration error between the image pair [30]. The former is intuitive and efficient while
 244 the latter perceives the error directly at the cost of calculating registration error between
 245 all available image pairs in advance. Considering the association between the number of
 246 matched features and the registration error, we creatively set the weight of an edge in
 247 the graph as follows:

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

248 where $M_{i,j}$ denotes the total number of matches between \mathbf{I}_i and \mathbf{I}_j , and ε is a constant
 249 for regularization, which is set as 50 in our experiments. This weight setting equation,
 250 which describes the contribution of matched features to the registration accuracy, has the
 251 advantages of both efficiency and effect.

252 Based on the weighted graph, the optimal reference image selection problem is
 253 formulated as finding a node with the least total weight of the shortest paths to all
 254 the other nodes, which can be solved by the Floyd Warshalls all-pairs shortest path

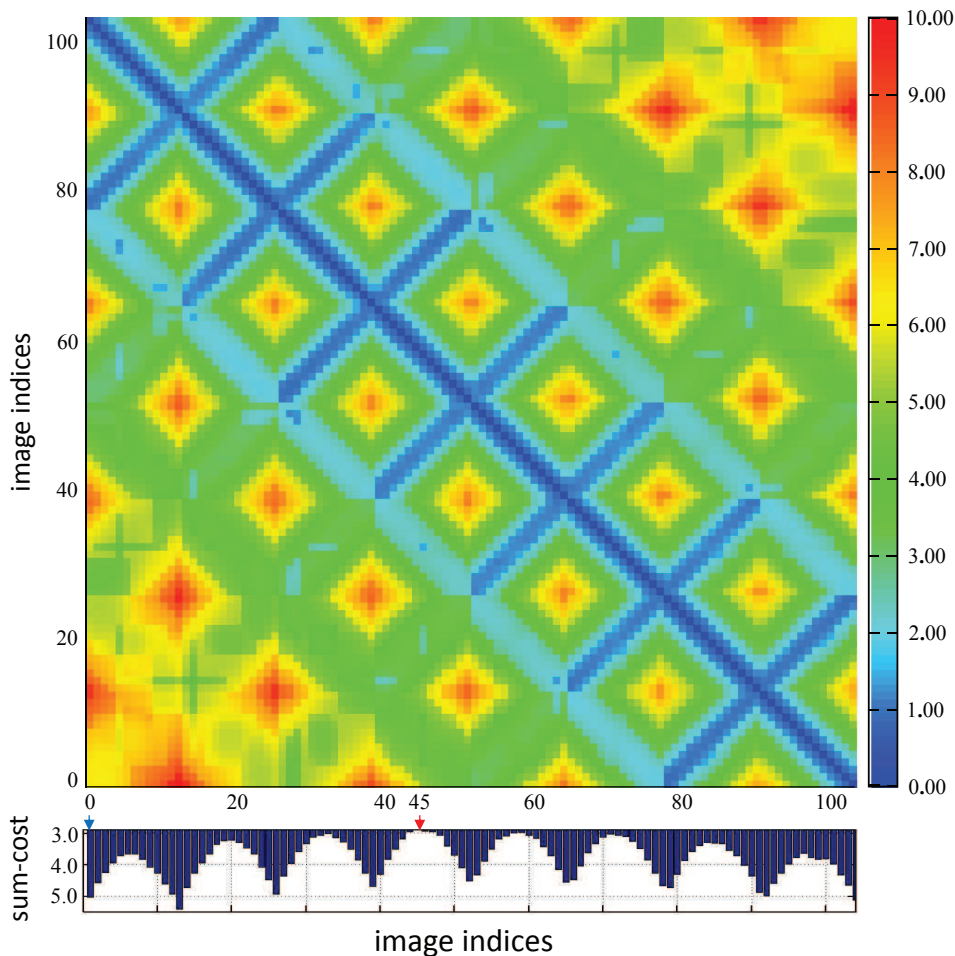


Figure 3: The cost matrix of all-pairs shortest path calculated from a sequence of 104 images. Below is attached the bar chart depicting the mean cumulative cost of each column of the cost matrix. As labeled with a red arrow in the bottom indices, the 45-th column of the cost matrix has the minimum total cost with the mean cumulative cost of 3.04. That’s to say, the 45-th image is the optimal reference image. However, the conventional idea to naively select the first image as the reference image gives an much higher mean cumulative cost of 5.23, as labeled with a blue arrow in the bottom indices.

255 algorithm [33, 34]. The dynamic programming strategy is applied in this algorithm with
 256 the computation complexity of $O(3)$, so it is more efficient than running n times of a
 257 single source shortest path algorithm. With this algorithm, all shortest paths from a
 258 node to any other node can be obtained. When there are n images in a sequence, we
 259 build a $n \times n$ size symmetric cost matrix \mathbf{W} where each element records the cost of the
 260 shortest path between two images. After running this algorithm, the cost of the shortest
 261 path from \mathbf{I}_i to \mathbf{I}_j is saved in $\mathbf{W}(i, j)$. Therefore, the i -th row or column of matrix \mathbf{W}

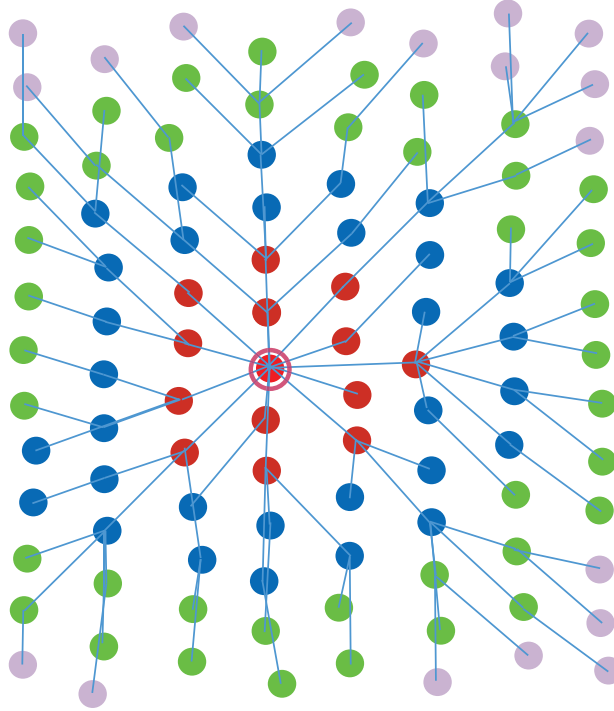


Figure 4: The spinning tree of the graph with the optimal reference image as the root node (marked with red ring). Nodes in different levels of the tree are marked with different colors, and the blue lines link each node and its parent node, which imply the shortest paths from all the other images to the reference image.

262 indicates the cost of every shortest path from other images to \mathbf{I}_i . On this occasion, the
 263 accumulated cost of each column in the cost matrix \mathbf{W} can be calculated and the column
 264 with the minimum accumulated cost is selected as the reference image. To demonstrate
 265 the procedure, the cost matrix \mathbf{W} of a sequence of 104 images is visualized in Figure 3,
 266 and the 45-th column with the minimum total cost is labeled with the red arrow in the
 267 bottom indices. Specially, the conventional strategy of selecting the first image as the
 268 reference image is also highlighted as a comparison. Considering the amount of images,
 269 the gap of the mean cumulative costs between the two strategies can make a big difference
 270 to the mosaicking result.

271 Actually, each row in \mathbf{W} corresponds a spinning tree with the node of this row as the
 272 root node, which describes the hierarchical relationship of the image nodes. With the
 273 selected reference image, the spinning tree of the image sequence described in Figure 3,
 274 is displayed in Figure 4. The spinning tree indicates the direct reference image of each

275 image (parent node in graphic terms), which determines the aligning order of images in
276 the following global alignment.

277 *2.3. Globally Consistent Alignment*

278 In general, both the locally aligning accuracy and the global consistency are two basic
279 factors determining the quality of mosaicking result. Under a strict transformation model,
280 these two factors can be guaranteed in a coherent way, where the higher aligning precision
281 contributes on the better global consistency. However, in most practical applications, the
282 observing scenes of pseudo-planes make the frequently-used homography model just an
283 approximate transformation between images. In this case, the aligning model of higher
284 degrees of freedom (DoF) usually makes more accurate alignment but suffers more severe
285 perspective distortion meanwhile, and vice versa. Therefore, we have to deal with these
286 two factors in a trade-off way. To keep the optimal balance between them, the model
287 with a relatively low DoF is employed to make the initial alignment of images robustly,
288 which will be refined with a higher DoF to improve the aligning precision under the
289 anti-perspective constraint.

290 *2.3.1. Robust Alignment by Affine Model*

291 For a robustly initial alignment, we would rather to use the affine model which
292 compromises between the 2D rigid transformation and the homography transformation.
293 On the one hand, the approximately coplanar constraint of images is partly implied in
294 the six-parameter affine model which can suppress severe perspective distortion to some
295 extent, on the other hand, the affine transformation is able to provide a qualified initial
296 solution for the following homography refinement.

297 According to the spinning tree mentioned in Section 2.2, the sequential images are
298 aligned group by group in the order of breadth-first search, which can decrease the
299 accumulation error of alignment compared to the way of one by one. In this paper, when
300 aligning a new group of images to the reference frame, the overlapping relations between
301 all the previously aligned images and the newly added images, and the overlapping
302 relations between intra-group will be jointly used in the optimization framework. Let
303 $\mathcal{I} = \{\mathbf{I}_i\}_{i=1}^s$ be the set of previously aligned images. The affine transformation set

304 $\mathcal{A} = \{\mathbf{A}_i\}_{i=s+1}^{s+m}$ of the newly added image group $\mathcal{G} = \{\mathbf{I}_i\}_{i=s+1}^{s+m}$ for alignment will be
 305 optimized by minimizing the combination of two cost functions as below :

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

306 where the first energy term $E_1(\mathcal{A}|\mathcal{I}, \mathcal{G})$ is related to the overlapping relations between \mathcal{I}
 307 and \mathcal{G} as follows:

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

308 and the second energy term $E_2(\mathcal{A}|\mathcal{G})$ is related to the overlapping relations in \mathcal{G} as follows:

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

309 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
 310 Section 2.

311 As a group of linear equations, Eq. (3) can be solved fast by the Singular Value
 312 Decomposition (SVD) method. Note that both the epipolar constraint and the
 313 appropriately homography constraint are employed to remove outliers in the SURF
 314 points matching algorithm. In fact, we also normalize the coordinates of matched points
 315 according to the method proposed in [35], in order to increase the numerical stability
 316 by improving the condition number of the coefficient matrix. What's more, the robust
 317 estimator MLESAC [36] is used to exclude outliers for affine estimation because it is
 318 beneficial for the image mosaicking of quasi-planar scenes.

319 2.3.2. Model Refinement under Anti-Perspective Constraint

320 The affine models recovered by groups are mainly used to achieve the robust initial
 321 alignment, which guarantee the mosaicking result against the perspective distortion well.
 322 However, the aligning precision needs to be further improved due to that the DoF of
 323 the aligning model is limited and no global optimization is performed. To improve the
 324 aligning accuracy to some extent but not to induce the perspective distortion, the energy
 325 function should allow to transit the affine model to the homography model with a higher
 326 DoF under some reasonable constraint. In fact, such constraint has been implied in
 327 the affine model which has the anti-perspective property relative to the homography

328 model. So, the deviation between the optimal homography transformation and the
 329 initially estimated affine transformation is set as a regularization term in the proposed
 330 optimization framework.

331 As the images are aligned by groups, the affine models of all the images $\mathcal{I} = \{\mathbf{I}_i\}_{i=1}^n$
 332 can be obtained, denoted as $\mathcal{A} = \{\mathbf{A}_i\}_{i=1}^n$, which are used as the initial parameters
 333 for the homography model in the final global optimization. The homography models
 334 $\mathcal{H} = \{\mathbf{H}_i\}_{i=1}^n$ with respect to the reference frame will be optimized in the energy
 335 function composed of two mutually contrary terms. The data term set for minimizing
 336 the sum of squares of the feature registration errors between images is denoted as:

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

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

341 Another optimization objective is to keep the global consistency by suppressing the
 342 accumulation of the perspective distortions which may emerge in the transition from
 343 the affine model to the homography model. The regularization term from the idea that
 344 the optimal homography transformation should be close to the initially estimated affine
 345 transformation, is expressed as the displacements of the warped features from their initial
 346 positions as follows:

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

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

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

351 where λ is the weight coefficient used for balancing the two terms E_d and E_r , which should
 352 be set to an appropriate small value since the constraint isn't a strict one. Theoretically,

353 a bigger value of λ strengthens the global consistency while decreases the accuracy of the
354 local alignment. We set its value from 0.01 to 0.05 in all our experiments. As a typical
355 non-linear least squares problem, Eq. (8) can be solved by the Levenberg-Marquardt (LM)
356 algorithm. However, considering the specialty of this problem, we employ the sparse LM
357 algorithm [37] to save memory and to speed up the computation, which is stated detailedly
358 in Appendix A.

359 3. Experimental Results

360 To make a comprehensive study of our approach, three groups of experiments were
361 conducted, including the evaluation on the topology estimation, the evaluation on the
362 selection of initial model, and the comprehensive evaluation on the mosaicking results.
363 Two sets of representative aerial images acquired by different flight platforms and over
364 different landforms, respectively, were used as the experimental dataset. The first dataset,
365 consisting of 744 images from 24 sequentially ordered strips, was captured at a flight height
366 of about 780 meters over an urban area. The original images, with a forward overlapping
367 rate of about 60% were down-sampled to the size of 1000×642 in our experiments. The
368 second dataset, consisting of 130 images with the down-sampling size of 800×533 , was
369 captured by an unmanned aerial vehicle (UAV) with a forward overlapping rate of about
370 70%, which observes a suburb area containing mountains.

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

374 3.1. Evaluation on Topology Estimation

375 In this section, the topology estimation module of the proposed approach was
376 compared with the classic all-against-all strategy and the state-of-the-art algorithm
377 implemented according to [28] (we name it as Fast-Topology hereafter). The comparisons
378 were performed on the estimated topology of the aforementioned two datasets. To
379 test our approach diversely, the aerial image sequence and the UAV image sequence
380 were respectively processed in two different modes for topology estimation: the time-
381 consecutive mode and the unordered mode. As a robust but exhaustive strategy,

382 matching all-against-all was always used for comparison in topology estimation, the
383 detected overlapping pairs by which can be regarded as the ground truth. Moreover,
384 the successfully matched image pairs and the total matching attempts are combined to
385 evaluate the topology estimation results as the quantitative metrics.

386 The topology estimation results of the two datasets are summarized in Table 2 and
387 Table 3, respectively, where the first column lists the tested methods, the second column
388 corresponds to the numbers of successfully matched image pairs, the third column contains
389 the total numbers of matching attempts, the last two columns give the percentages of the
390 second and third columns with respect to the all-against-all strategy. As the tables show,
391 both our approach and Fast-Topology [28] almost recovered the complete topology as the
392 all-against-all strategy did, but with a much less amount of image matching attempts.
393 Although there are some omissions with respect to all-against-all, the major overlapping
394 relations had been detected successfully in our approach, which can be observed in the
395 topological graph depicted in Figure 5(a) and Figure 6(a), respectively. It implies that
396 most of the undetected overlapping pairs probably share very small overlapping areas
397 and make little difference to the mosaicking results. Compared to Fast-Topology [28],
398 our approach has roughly the same recall rates but less total matching attempts, which
399 benefits from two key strategies used in the potential overlapping pairs detection. The
400 one is the selection the temporary reference image, which is determined by applying the
401 strategy detailed in Section 2.2 on the main chain, instead of setting the first image
402 simply like Fast-Topology. The other is that the position of the newly added image is
403 simultaneously adapted along with the potential overlapping relations being detected,
404 which improves the alignment accuracy and so does the efficiency. Differently with ours,
405 the procedures of detecting the potential overlapping pairs and adapting alignment of
406 images with the detecting results are divided into two independent steps in Fast-Topology.
407 Therefore, it inevitably introduced many unnecessary matching attempts because of the
408 inaccurate alignment in the first few iterations though it can find most of the existing
409 overlapping relations after several iterations.

410 As mentioned in Section 2.2, the estimated topology is used to search for the optimal
411 reference image, by the way of which the images are organized as a spinning tree implying

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

412 the aligning order for the global alignment. Here, the spinning trees with the reference
413 image as the root node, are expressed by a group of red edges of the topological graph in
414 Figure 5(b) and Figure 6(b), corresponding to the first and second datasets, respectively.
415 It's easy to find that the selected reference images can always locate in the central
416 part geometrically, no matter of the ruled aerial data or the strip shaped UAV data.
417 Noticeably, the layouts of the image centroids recovered via the final global alignment,
418 depicted in Figure 5(b) and Figure 6(b), are more neat (accurate) than those depicted
419 in Figure 5(a) and Figure 6(a), respectively. This is reasonable, because the global
420 alignment for compositing a good mosaic includes both the topology analysis and the
421 global optimization while the topology estimation just aims at finding the topology in an
422 efficient way.

423 3.2. Evaluation on Initial Model Selection

424 In the period of recovering initial alignment described in Section 2.3.1, the selection
425 of the transformation model among *rigid*, *affine* and *homography* models can make
426 differences to the final mosaicking result. To amplify the influence of error factors, we

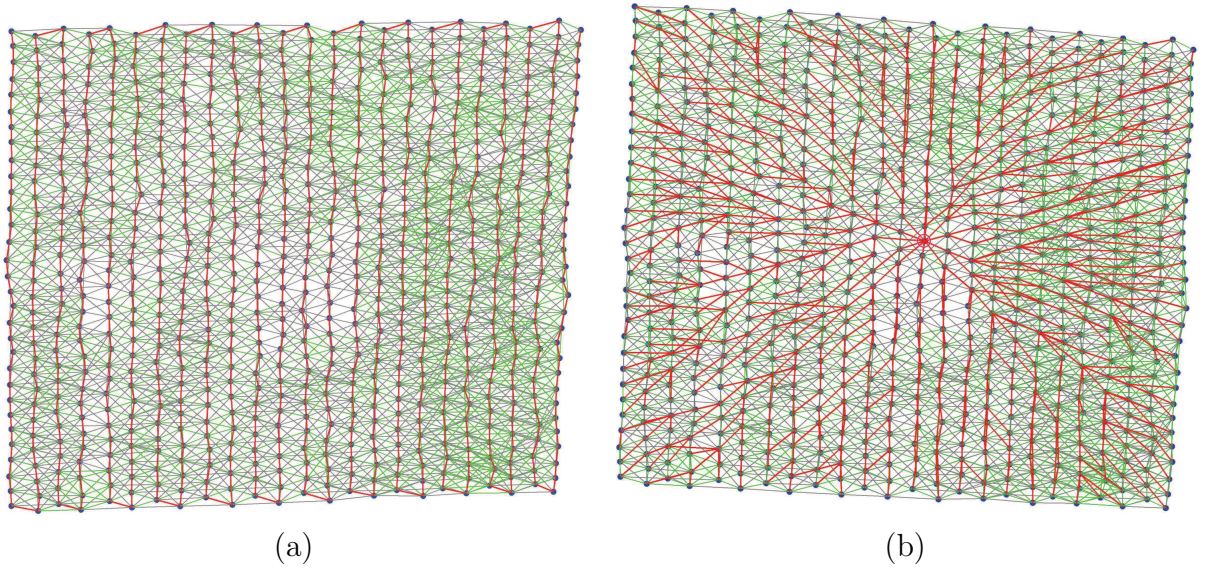


Figure 5: The estimated topology of the first dataset (744 images) highlighted for different aims: (a) The estimated topology with the prior main chain in the time-consecutive order marked with red edges; (b) The spinning tree generated by searching for the optimal reference image (the node with red ring), marked with red edges on the estimated topological graph. Different from (a), the geometric positions in (b) were recovered by the final global alignment. The edges in green and gray indicate the numbers of matched features between image pairs more and less than 100, respectively.

Table 4: 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

427 specially selected a strip-shaped aerial image subset and a block UAV image subset from
 428 the first dataset and the second one, respectively, and the image on the end was set as
 429 the reference image. The comparative analyses were made on both alignment precision
 430 and global consistency, where the numerical results are shown in Table 4 while the global
 431 consistency can be judged via the visual results shown in Figure 7.

432 As for the strip aerial images, the homography model employed as the initial model
 433 has the best alignment precision, but suffers severe an accumulation of the perspective

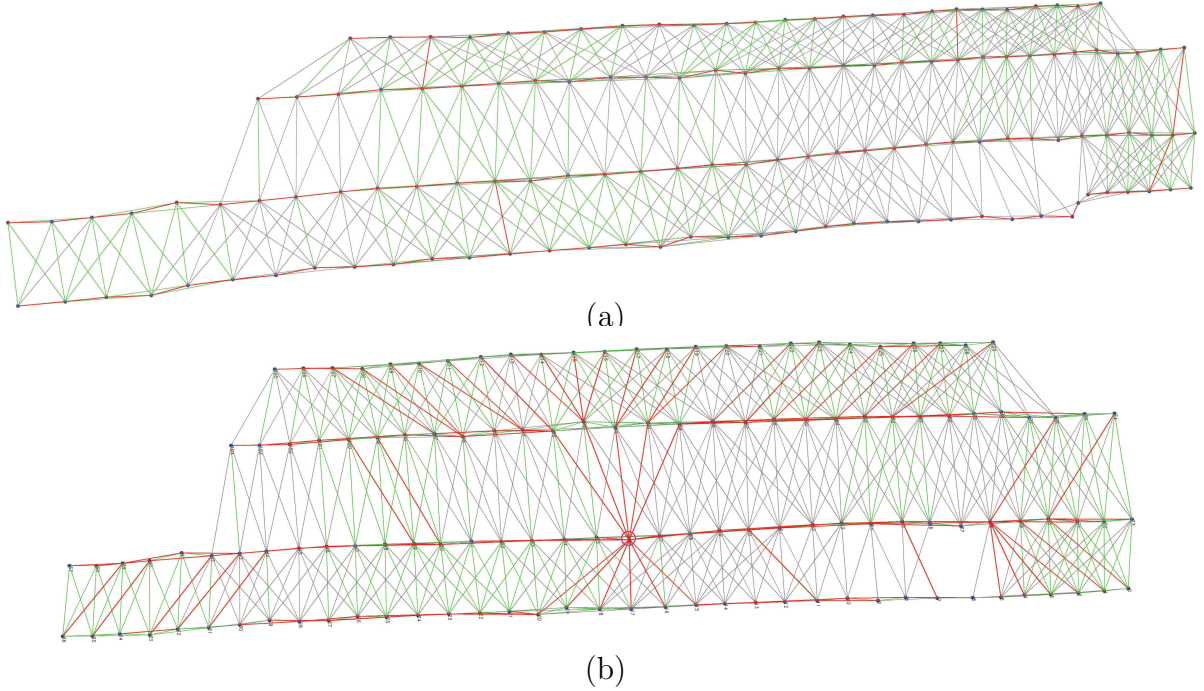


Figure 6: The estimated topology of the second dataset (130 images) highlighted for different aims: (a) The estimated topology with the main chain, linked by the proposed iterative scheme, labeled in red; (b) The spinning tree generated by searching for the optimal reference image (the node with red ring), marked with red edges on the estimated topological graph. Different from (a), the geometric positions in (b) were recovered by the final global alignment. The edges in green and gray indicate the numbers of matched features between image pairs more and less than 100, respectively.

434 distortions meanwhile due to that it has the highest DoF for alignment. However,
 435 the mosaicking result based on the rigid transformation as the initial model shows a
 436 bending tendency with the lowest accuracy although it doesn't induce a severe perspective
 437 distortion. This is because the rigid model of 3 free parameters just allows the image
 438 translation and rotation, which is not enough to describe the truly geometric relations
 439 between images and prone to resulting in the accumulation of rotation or translation.
 440 Compromising between them, the affine model with a moderate DoF, has made a good
 441 balance between the aligning accuracy and the global consistency, which gives the most
 442 visually satisfactory mosaicking result.

443 Because of the low flight altitude, the comparatively large-depth-difference ground
 444 greatly decreases the aligning precision for the UAV image sequence. In this case, the
 445 perspective distortion is still noticeable for the mosaicking result of the homography

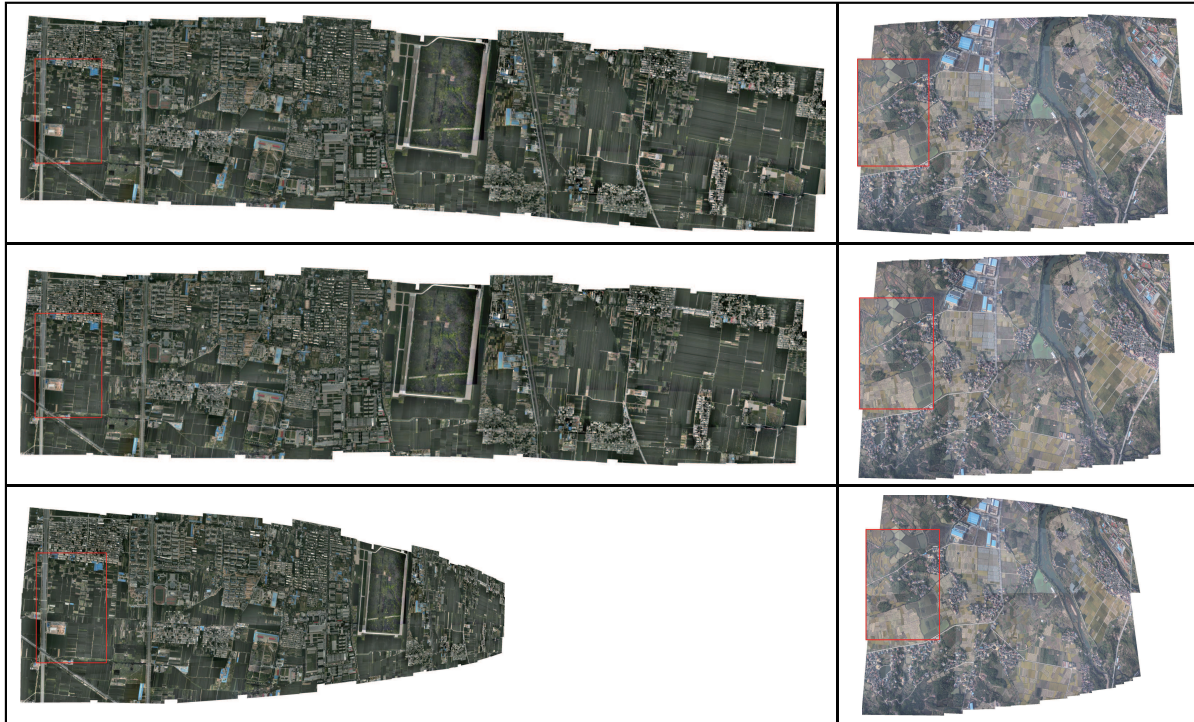


Figure 7: The thumbnails of the mosaicking results on the aerial images (Left) and the UAV images (Right) where the rigid model in the first row, the affine model in the second row, and the homography model in the last row were chosen for initial alignment, respectively. Notice that the reference image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

446 model, even though the images were taken from a small block area. Differently, the rigid
 447 model achieves an as good visual result as the affine model does for this UAV image
 448 sequence, though its aligning precision is a little inferior to that of the affine model.
 449 Conclusively, the affine model has the best comprehensive property to provide a robust
 450 initial alignment, so it's the most reasonable choice of the initial aligning model in our
 451 approach.

452 3.3. Comprehensive Evaluation on Mosaicking Results

453 The final mosaicking results of our approach were evaluated in both qualitative and
 454 quantitative forms. Firstly, we compared the mosaicked images generated by our approach
 455 with those created by a commercial software named PTGui ¹ on visual effects. Since
 456 aiming at comparing the alignment results only, the following seamline detection and

¹<http://www.ptgui.com/>

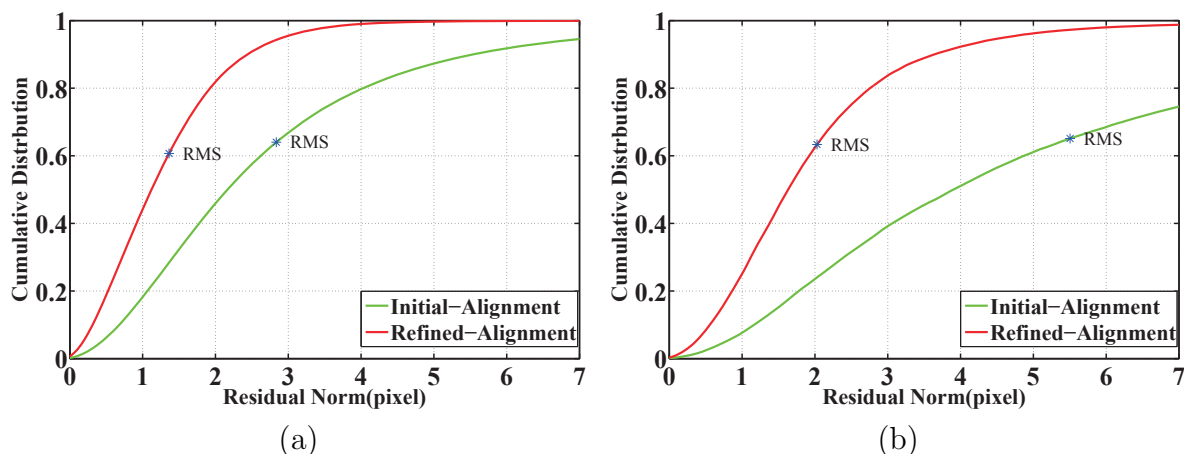


Figure 8: Cumulative probability distributions of the residual error norms with and without the global refinement performed in our approach: (a) the error analysis for the first dataset; (b) the error analysis for the second dataset. The green curves depict the aligning error of our approach with only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

457 tonal correction were skipped in PTGui and our image stacking order was made consistent
 458 with that of PTGui. The comparative results of the first and the second datasets are
 459 illustrated in Figure 9 and Figure 10, respectively.

460 From the mosaics shown in Figure 9, the two mosaics have similar visual effects as a
 461 whole, both of which take on a pretty good global consistency. However, when it comes to
 462 the local aligning accuracy, our approach has an obvious superiority over PTGui, which
 463 can be observed from some enlarged regions listed in the right column of Figure 9. As for
 464 the UAV data, the large-depth-difference ground makes the assumption of planarity of
 465 the scene weaker, which increases the difficulty to keep the global consistency. A slightly
 466 down-scale tendency in the left part can be found in the mosaicking result of our approach
 467 in Figure 10(a). Since some strong constraints were employed for keeping the scale of
 468 each image consistent, the mosaicking result of PTGui nearly suffered no perspective
 469 distortions, but in the mean while, the alignment precision was destroyed greatly. For
 470 a detail comparison, a serial of enlarged typical regions are listed in the middle line of
 471 Figure 10, which illustrate the excellent performance of our approach in the aspect of
 472 aligning accuracy.

473 Without precision analysis in PTGui, the quantitative evaluation of our approach

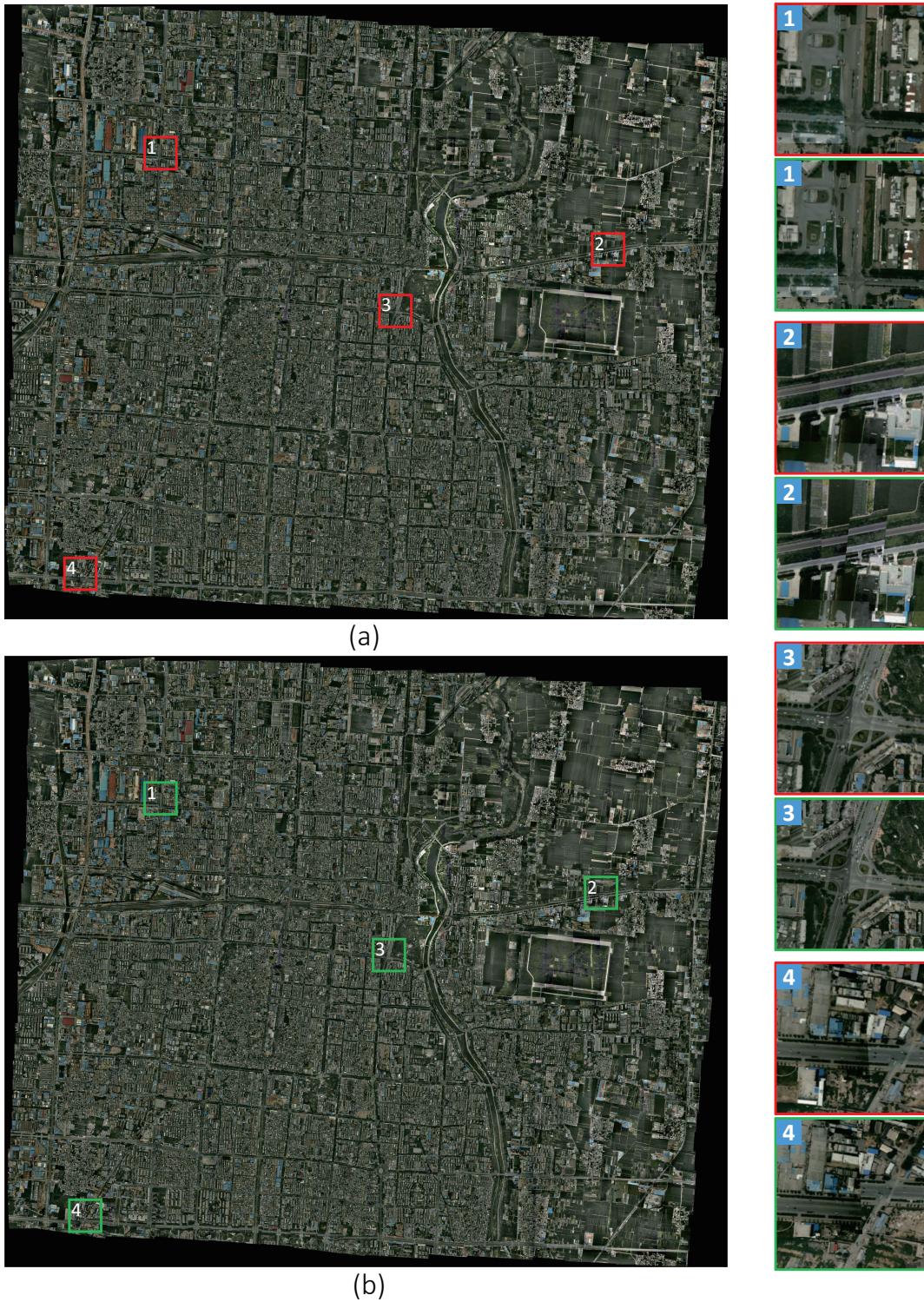


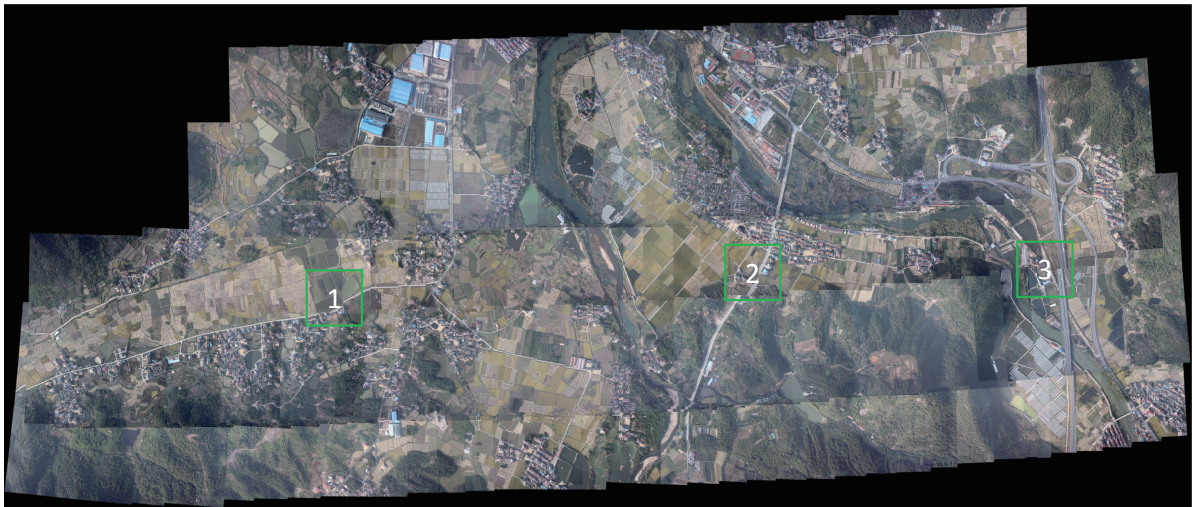
Figure 9: The mosaics composited from the first dataset (744 images) by: (a) our approach and (b) PTGui, respectively. Several typical regions grabbed from the mosaics are enlarged in pairs in the right column.



(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).

474 was performed in two aspects. As an alignment precision, the registration error of our
 475 approach, running with the initial alignment only and the full set of global alignment,
 476 respectively, were compared in the form of cumulative probability distribution, as

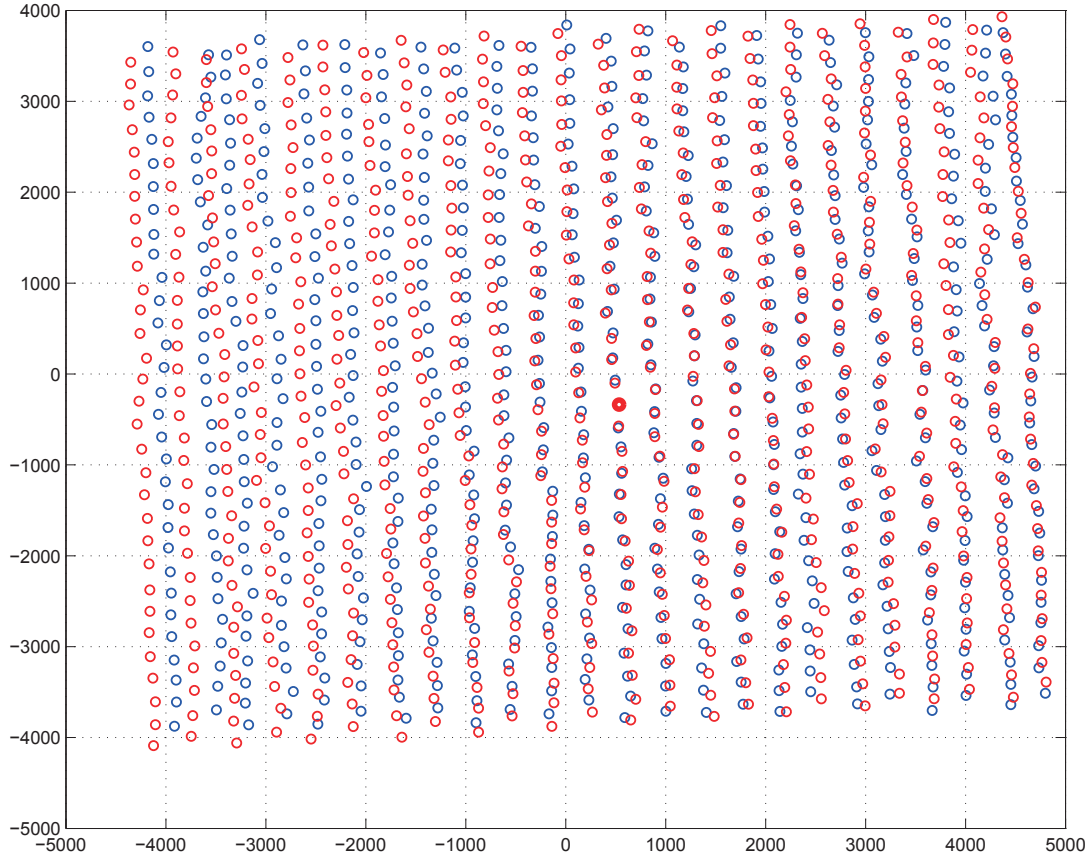


Figure 11: Distributions of image centroids on the mosaic computed by two different approaches. The red circles are the centroids recovered by the proposed approach, and the blue ones represent the result of the pose-based approach. The solid red circle stands for the centroid of the reference image,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

477 displayed in Figure 8. From the comparisons, it's easy to find that the aligning precision
 478 increases a lot with the help of the homography refinement, while the global consistency
 479 is not affected during the transition from the affine model to the homography one, as can
 480 be observed in Figure 9(a) and Figure 10(a). This is what we try to achieve, namely to
 481 keep an optimal balance between the alignment accuracy and the global consistency.

482 Moreover, the available poses of the first dataset, which were recovered by the rigid
 483 block adjustment of photogrammetry field, were used to calculate the homography models
 484 according to the formula in [4] under the assumption of the ground being a plane.
 485 Considering the pseudo-planarity of the ground scene, they are not accurate enough
 486 to be used as the ground truth, but they are qualified to evaluate the global consistency

487 as a reference, since the pose parameters can be regarded as no accumulation error. The
488 recovered image centroids of the first dataset obtained by our approach and the reference
489 models, are illustrated in Figure 11. It shows that the two groups of centroids have a
490 similar distribution form but there are also some displacements between corresponding
491 centroids which average at 5.16 pixels. Obviously, the displacements in the right part are
492 much smaller than those in the left part, because the right part has more dense and strong
493 topological relationships suppressing accumulation errors. In fact, as an image mosaicking
494 approach based on the 2D feature registration, the recovered geometric positions are
495 accurate enough to keep the global consistency of a mosaic, which emphasizes more on
496 the visual effects than the geometric measurements. What’s more, because of no image
497 registration based optimization performed, the pose-based approach obtains a terrible
498 image aligning accuracy as the RMS error of 103.9 pixels, which is much inferior to that
499 of 1.36 pixels in our approach. Therefore, our approach has a good property of alignment
500 accuracy and global consistency in the terms of image mosaicking.

501 4. Conclusion and Future Works

502 In this paper, a topology analysis based generic framework was proposed for
503 mosaicking sequential images of an approximately planar scene, which is composed
504 of three steps : topology estimation, reference image selection, and global alignment.
505 Specifically, it’s adapted to both ordered and unordered image sequences. To estimate
506 the topology robustly, we perform the image location and the potential overlapping
507 pairs detection in a collaborative way, which results in that our approach for topology
508 estimation significantly outperformed the state-of-the-art method in the aspect of
509 efficiency. Based on the topological graph, the optimal reference image is found by graph
510 analysis and all the images are organized as a spinning tree which gives the reference
511 relationships for each image. With the result of topological analysis, we propose a global
512 alignment strategy of allowing the continuous transition between the affine model and the
513 homography one according to the energy definition, which can keep the optimal balance
514 between the global consistency and the aligning accuracy adaptively. The proposed
515 framework was tested with several datasets and the experimental results illustrate the

516 superiority of our approaches. However, as stated in this paper, the global consistency
517 and the alignment accuracy need to be treated in a trade-off way in the case of pseudo-
518 plane. Therefore, the ideal following process of this would be optimal seamline selection,
519 which removes the residual parallax by crossing areas with less misalignment. This is
520 meaningful for compositing a mosaicked image of high quality, and it will be studied in
521 the future work.

522 Acknowledgment

523 This work was partially supported by the National Natural Science Foundation of
524 China (Project No. 41571436), the Hubei Province Science and Technology Support
525 Program, China (Project No. 2015BAA027), the National Natural Science Foundation of
526 China (Project No. 41271431), and the Jiangsu Province Science and Technology Support
527 Program, China (Project No. BE2014866).

528 Appendix A. Optimization Derivation for Model Refinement Under Anti- 529 Perspective Constraint

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

534 Here, we define the parameter vector of the homography matrix \mathbf{H}_i as $\theta_i =$
535 $[h_1^i, h_2^i, h_3^i, h_4^i, h_5^i, h_6^i, h_7^i, h_8^i]^\top$, $i \in [1, n]$, and the initial value of θ_i is defined as $\bar{\theta}_i =$
536 $[\bar{h}_1^i, \bar{h}_2^i, \bar{h}_3^i, \bar{h}_4^i, \bar{h}_5^i, \bar{h}_6^i, \bar{h}_7^i, \bar{h}_8^i]^\top$. Taking a pair of matching points $\{\varpi(\mathbf{x}_{ij}^k) = (x, y), \varpi(\mathbf{x}_{ji}^k) =$
537 $(x', y')\}$ from \mathbf{I}_i and \mathbf{I}_j for example, Eq. (8) can be written as :

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

538 where $[x_0, y_0]^\top = \varpi(\mathbf{A}_1 \mathbf{x}_{ij}^k)$ and $[x'_0, y'_0]^\top = \varpi(\mathbf{A}_j \mathbf{x}_{ji}^k)$, are the constant terms which can
539 be calculated in advance. Eq. (A.1) is expanded in the form of the first-order Taylor
540 series as:

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

541 where \bar{f}_k is the values of f_k when substituting $\bar{\theta}_i$ and $\bar{\theta}_j$ into Eq. (A.1). $d\theta_i =$
542 $[dh_1^i, dh_2^i, dh_3^i, dh_4^i, dh_5^i, dh_6^i, dh_7^i, dh_8^i]^\top$ represents the delta value of $\theta_i, i \in [1, n]$. The
543 partial derivatives of functions f_k with respect to θ_i and θ_j are listed as below:

$$\left\{ \begin{array}{l} \frac{\partial f_k}{\partial h_1^i} = \frac{K_1 x}{\bar{h}_7^i x + \bar{h}_8^i y + 1}, \quad \frac{\partial f_k}{\partial h_2^i} = \frac{K_1 y}{\bar{h}_7^i x + \bar{h}_8^i y + 1}, \quad \frac{\partial f_k}{\partial h_3^i} = \frac{K_1}{\bar{h}_7^i x + \bar{h}_8^i y + 1}, \\ \frac{\partial f_k}{\partial h_4^i} = \frac{K_2 x}{\bar{h}_7^i x + \bar{h}_8^i y + 1}, \quad \frac{\partial f_k}{\partial h_5^i} = \frac{K_2 y}{\bar{h}_7^i x + \bar{h}_8^i y + 1}, \quad \frac{\partial f_k}{\partial h_6^i} = \frac{K_2}{\bar{h}_7^i x + \bar{h}_8^i y + 1}, \\ \frac{\partial f_k}{\partial h_7^i} = \frac{-K_1(\bar{h}_1^i x + \bar{h}_2^i y + \bar{h}_3^i)x}{(\bar{h}_7^i x + \bar{h}_8^i y + 1)^2} + \frac{-K_2(\bar{h}_3^i x + \bar{h}_4^i y + \bar{h}_5^i)x}{(\bar{h}_7^i x + \bar{h}_8^i y + 1)^2}, \\ \frac{\partial f_k}{\partial h_8^i} = \frac{-K_1(\bar{h}_1^i x + \bar{h}_2^i y + \bar{h}_3^i)y}{(\bar{h}_7^i x + \bar{h}_8^i y + 1)^2} + \frac{-K_2(\bar{h}_3^i x + \bar{h}_4^i y + \bar{h}_5^i)y}{(\bar{h}_7^i x + \bar{h}_8^i y + 1)^2}, \\ \frac{\partial f_k}{\partial h_1^j} = \frac{K_3 x}{\bar{h}_7^j x + \bar{h}_8^j y + 1}, \quad \frac{\partial f_k}{\partial h_2^j} = \frac{K_3 y}{\bar{h}_7^j x + \bar{h}_8^j y + 1}, \quad \frac{\partial f_k}{\partial h_3^j} = \frac{K_3}{\bar{h}_7^j x + \bar{h}_8^j y + 1}, \\ \frac{\partial f_k}{\partial h_4^j} = \frac{K_4 x}{\bar{h}_7^j x + \bar{h}_8^j y + 1}, \quad \frac{\partial f_k}{\partial h_5^j} = \frac{K_4 y}{\bar{h}_7^j x + \bar{h}_8^j y + 1}, \quad \frac{\partial f_k}{\partial h_6^j} = \frac{K_4}{\bar{h}_7^j x + \bar{h}_8^j y + 1}, \\ \frac{\partial f_k}{\partial h_7^j} = \frac{-K_3(\bar{h}_1^j x + \bar{h}_2^j y + \bar{h}_3^j)x}{(\bar{h}_7^j x + \bar{h}_8^j y + 1)^2} + \frac{-K_4(\bar{h}_3^j x + \bar{h}_4^j y + \bar{h}_5^j)x}{(\bar{h}_7^j x + \bar{h}_8^j y + 1)^2}, \\ \frac{\partial f_k}{\partial h_8^j} = \frac{-K_3(\bar{h}_1^j x + \bar{h}_2^j y + \bar{h}_3^j)y}{(\bar{h}_7^j x + \bar{h}_8^j y + 1)^2} + \frac{-K_4(\bar{h}_3^j x + \bar{h}_4^j y + \bar{h}_5^j)y}{(\bar{h}_7^j x + \bar{h}_8^j y + 1)^2}, \end{array} \right.$$

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

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

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

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

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

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

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

552 As can be seen, a pair of matching points from two images provides an equation with
 553 16 unknown parameters. Supposing that n images have m pairs of overlapping relations
 554 and there are s matching points of each image pair in average, then we obtain a Jacobi
 555 matrix with the size of $m \times s$ rows and $8 \times m$ columns and a constant matrix with the size
 556 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
 557 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
 558 following equation:

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

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

564 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 $\mathbf{TL}_{8n \times 1} =$
 565 $\mathbf{J}_{ms \times 8n}^\top \mathbf{L}_{ms \times 1}$ have been obtained. So, to reduce the required memory space and the

566 computation time, we manage to compute $\mathbf{TJ}_{8n \times 8n}$ and $\mathbf{TL}_{8n \times 1}$ by adding up the matrix
567 $\mathbf{J}^{i\top} \mathbf{J}^i$ and the matrix $\mathbf{J}^{i\top} \mathbf{L}^i$ calculated from each pair of matching points, instead of
568 building the large \mathbf{J} and \mathbf{L} beforehand. The improved computation formula is defined as:

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

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

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

570 What's more, considering the sparsity of \mathbf{J}^i , the computation of the matrix multiplication
571 in Eq. (A.5) can be improved further in the complexity of both time and space.

572 References

- 573 [1] E. Zagrouba, W. Barhoumi, S. Amri, An efficient image-mosaicing method based on
574 multifeature matching, *Machine Vision and Applications* 20 (3) (2009) 139–162.
- 575 [2] J. Chen, H. Feng, K. Pan, Z. Xu, Q. Li, An optimization method for registration and
576 mosaicking of remote sensing images, *Optik - International Journal for Light and Electron*
577 *Optics* 125 (2) (2014) 697–703.
- 578 [3] B. Zitova, J. Flusser, Image registration methods: a survey, *Image and Vision Computing*
579 21 (11) (2003) 977–1000.
- 580 [4] L. Kang, L. Wu, Y. Wei, B. Yang, H. Song, A highly accurate dense approach for
581 homography estimation using modified differential evolution, *Engineering Applications of*
582 *Artificial Intelligence* 31 (4) (2014) 68–77.
- 583 [5] Z. Wang, Y. Chen, Z. Zhu, W. Zhao, An automatic panoramic image mosaic method based
584 on graph model, *Multimedia Tools and Applications* 75 (5) (2015) 2725–2740.
- 585 [6] N. R. Gracias, S. Van Der Zwaan, A. Bernardino, S.-V. Jos, Mosaic-based navigation
586 for autonomous underwater vehicles, *IEEE Journal of Oceanic Engineering* 28 (4) (2003)
587 609–624.

- 588 [7] Y. Xu, J. Ou, H. He, X. Zhang, J. Mills, Mosaicking of unmanned aerial vehicle imagery
589 in the absence of camera poses, *2016 8 (3) (2016) 204*.
- 590 [8] Y. He, R. Chung, Image mosaicking for polyhedral scene and in particular singly visible
591 surfaces, *Pattern Recognition 41 (3) (2008) 1200–1213*.
- 592 [9] J. Zaragoza, T.-J. Chin, M. Brown, D. Suter, As-projective-as-possible image stitching
593 with moving DLT, *IEEE Transactions on Pattern Analysis and Machine Intelligence 36 (7)*
594 *(2014) 1285–1298*.
- 595 [10] C.-H. Chang, Y. Sato, Y.-Y. Chuang, Shape-preserving half-projective warps for image
596 stitching, in: *IEEE Conference on Computer Vision and Pattern Recognition, 2014*, pp.
597 3254–3261.
- 598 [11] D. Patidar, A. Jain, Automatic image mosaicing: an approach based on FFT, *International*
599 *Journal of Scientific Engineering and Technology 1 (1) (2011) 01–04*.
- 600 [12] S. Ghannam, A. L. Abbott, Cross correlation versus mutual information for image
601 mosaicing, *International Journal of Advanced Computer Science and Applications 4 (11)*
602 *(2013) 94–102*.
- 603 [13] A. Elibol, R. Garcia, O. Delaunoy, N. Gracias, Efficient Topology Estimation for Large
604 Scale Optical Mapping, Springer, 2013, Ch. A New Global Alignment Method for Feature
605 Based Image Mosaicing, pp. 25–39.
- 606 [14] S. Ali, C. Daul, E. Galbrun, Guillemin, Anisotropic motion estimation on edge preserving
607 Riesz wavelets for robust video mosaicing, *Pattern Recognition 51 (2016) 425–442*.
- 608 [15] Y. Chen, J. Sun, G. Wang, Minimizing geometric distance by iterative linear optimization,
609 in: *IEEE International Conference on Pattern Recognition, Vol. 29, 2010*, pp. 1–4.
- 610 [16] W. Mou, H. Wang, G. Seet, L. Zhou, Robust homography estimation based on nonlinear
611 least squares optimization, *Mathematical Problems in Engineering 2014 (1) (2013) 372–377*.
- 612 [17] L. Zhang, Y. Li, J. Zhang, Y. Hu, Homography estimation in omnidirectional vision under
613 the L_∞ -norm, in: *IEEE International Conference on Robotics and Biomimetics, 2010*, pp.
614 1468–1473.

- 615 [18] B. Triggs, P. F. McLauchlan, R. I. Hartley, A. W. Fitzgibbon, *Vision Algorithms: Theory*
616 *and Practice, Lecture Notes in Computer Science*, Springer, 2000, Ch. Bundle adjustment
617 modern synthesis, pp. 298–372.
- 618 [19] K. Konolige, Sparse sparse bundle adjustment, in: *British Machine Vision Conference*,
619 2010, pp. 1–10.
- 620 [20] M. Li, D. Li, D. Fan, A study on automatic UAV image mosaic method for paroxysmal
621 disaster, in: *Proceedings of the International Society of Photogrammetry and Remote*
622 *Sensing Congress*, 2012.
- 623 [21] C. Xing, J. Wang, Y. Xu, A robust method for mosaicking sequence images obtained from
624 UAV, in: *International Conference on Information Engineering and Computer Science*
625 (ICIECS), 2010.
- 626 [22] F. Caballero, L. Merino, J. Ferruz, A. Ollero, Homography based Kalman filter for mosaic
627 building. applications to UAV position estimation, in: *IEEE International Conference on*
628 *Robotics and Automation*, 2007, pp. 2004–2009.
- 629 [23] A. Y. Taygun Kekec, M. Unel, A new approach to real-time mosaicing of aerial images,
630 *Robotics and Autonomous Systems* 62 (12) (2014) 1755–1767.
- 631 [24] E.-Y. Kang, I. Cohen, G. Medioni, A graph-based global registration for 2D mosaics, in:
632 *International Conference on Pattern Recognition*, Vol. 1, 2000, pp. 257–260.
- 633 [25] R. Marzotto, A. Fusiello, V. Murino, High resolution video mosaicing with global alignment,
634 in: *IEEE Conference on Computer Vision and Pattern Recognition*, Vol. 1, 2004, pp. 692–
635 698.
- 636 [26] A. Elibol, R. Garcia R Elibol A, Gracias Na, O. Delaunoy, N. Gracias, A new global
637 alignment method for feature based image mosaicing, *Advances in Visual Computing*,
638 *Lecture Notes in Computer Science* 5359 (7) (2008) 257–266.
- 639 [27] M. Xia, J. Yao, L. Li, X. Lu, Globally consistent alignment for mosaicking aerial images,
640 in: *IEEE International Conference on Image Processing*, 2015.

- 641 [28] A. Elibol, N. Gracias, R. Garcia, Fast topology estimation for image mosaicing using
642 adaptive information thresholding, *Robotics and Autonomous systems* 61 (2) (2013) 125–
643 136.
- 644 [29] R. Szeliski, Image alignment and stitching: A tutorial, *Foundations and Trends in*
645 *Computer Graphics and Vision* 2 (1) (2006) 1–104.
- 646 [30] T. E. Choe, I. Cohen, M. Lee, G. Medioni, Optimal global mosaic generation from retinal
647 images, in: *IEEE International Conference on Pattern Recognition*, Vol. 3, 2006, pp. 681–
648 684.
- 649 [31] B. Bollobas, *Modern graph theory*, Springer Science & Business Media, 2013.
- 650 [32] R. L. Graham, P. Hell, On the history of the minimum spanning tree problem, *Annals of*
651 *the History of Computing* 7 (1) (1985) 43–57.
- 652 [33] R. W. Floyd, Algorithm 97: shortest path, *Communications of the ACM* 5 (6) (1962)
653 345–345.
- 654 [34] T. H. Cormen, *Introduction to algorithms*, MIT press, 2009.
- 655 [35] R. I. Hartley, In defense of the eight-point algorithm, *IEEE Transactions on Pattern*
656 *Analysis and Machine Intelligence* 19 (6) (1997) 580–593.
- 657 [36] P. Torr, A. Zisserman, MLESAC: A new robust estimator with application to estimating
658 image geometry, *Computer Vision and Image Understanding* 78 (1) (2000) 138–156.
- 659 [37] M. I. A. Lourakis, Sparse non-linear least squares optimization for geometric vision, in:
660 *Computer Vision–ECCV 2010*, Vol. 6312, 2010, pp. 43–56.